

SymbolFit: Automatic Parametric Modeling with Symbolic Regression

Ho Fung Tsoi¹, Dylan Rankin¹, Cecile Caillol², Miles Cranmer³,
Sridhara Dasu⁴, Javier Duarte⁵, Philip Harris^{6,7}, Elliot Lipeles¹,
Vladimir Loncar^{6,8}

¹University of Pennsylvania, USA

²European Organization for Nuclear Research (CERN), Switzerland

³University of Cambridge, UK

⁴University of Wisconsin-Madison, USA

⁵University of California San Diego, USA

⁶Massachusetts Institute of Technology, USA

⁷Institute for Artificial Intelligence and Fundamental Interactions, USA

⁸Institute of Physics Belgrade, Serbia

E-mail: ho.fung.tsoi@cern.ch

Abstract. We introduce `SymbolFit`[†], a framework that automates parametric modeling by using symbolic regression to perform a machine-search for functions that fit the data, while simultaneously providing uncertainty estimates in a single run. Traditionally, constructing a parametric model to accurately describe binned data has been a manual and iterative process, requiring an adequate functional form to be determined before the fit can be performed. The main challenge arises when the appropriate functional forms cannot be derived from first principles, especially when there is no underlying true closed-form function for the distribution. In this work, we address this problem by utilizing symbolic regression, a machine learning technique that explores a vast space of candidate functions without needing a predefined functional form, treating the functional form itself as a trainable parameter. Our approach is demonstrated in data analysis applications in high-energy physics experiments at the CERN Large Hadron Collider (LHC). We demonstrate its effectiveness and efficiency using five real proton-proton collision datasets from new physics searches at the LHC, namely the background modeling in resonance searches for high-mass dijet, trijet, paired-dijet, diphoton, and dimuon events. We also validate the framework using several toy datasets with one and more variables.

[†]An API is available at <https://github.com/hftsoi/symbolfit>

1	Introduction	2
2	Challenges in traditional methods	3
2.1	Scenario 1: signal and background modeling	3
2.2	Scenario 2: derivation of smooth descriptions from binned data	4
2.3	An alternative method: Gaussian process regression	5
3	Proposed solution with symbolic regression	5
4	Method	6
5	Demonstrations	7
5.1	Toy Dataset 1 (1D) [signal modeling]	8
5.2	CMS dijet dataset (1D) [background modeling]	15
5.3	Toy dataset 3 (2D) [arbitrary shapes]	21
6	Summary	27
Appendix A	More examples	30
Appendix A.1	Toy dataset 2 (1D) [arbitrary shapes]	30
Appendix A.2	CMS high-mass diphoton dataset (1D) [background modeling]	34
Appendix A.3	CMS trijet dataset (1D) [background modeling]	39
Appendix A.4	CMS paired-dijet dataset (1D) [background modeling]	44
Appendix A.5	CMS high-mass dimuon dataset (1D) [background modeling]	49

1. Introduction

Traditional regression methods, such as polynomial regression, require specifying and fixing an adequate functional form before fitting the data. Identifying suitable functional forms for distributions with arbitrary shapes is often challenging and time-consuming, as, in most cases, these functions cannot be derived from first principles and must be determined through trial and error. Instead, symbolic regression (SR) is a more flexible and powerful technique that performs a *machine-search* for functions that best fit the data. In SR, the functional form itself is treated as a trainable parameter that is dynamically adjusted throughout the fitting process, eliminating the need to predefine an exact function—an empirical task that is often challenging [1].

Genetic programming [2] is a popular approach to SR [3–7]. In this approach, a function is represented as an expression tree, where the building blocks—mathematical operators, variables, and constants—are denoted as nodes, connected to represent their algebraic relations. Different functional forms are generated through the evolution these expression trees, where tree nodes are randomly selected and changed (mutation), and subtrees from different candidates are swapped to create new candidates (crossover), as illustrated in Fig. 1. As a result, the functional forms evolve during the fitting process, guiding the model toward convergence. Instead of predefining the final functional form, SR algorithms based on genetic programming need far less prior knowledge about the functions themselves. Only the constraints for constructing the expression trees need to be specified, such as the allowable mathematical operators (+, ×, /, pow, sin(·), exp(·), etc.). This flexibility eliminates the need to know the exact fitting function beforehand or to fine-tune one empirically.

In the following section, we identify two data analysis scenarios where parametric modeling is traditionally used in high-energy physics (HEP) experiments at the CERN Large Hadron Collider (LHC). We discuss the limitations of the current methods used in these contexts and how, without the need to predefine a functional form when fitting data, SR can help.

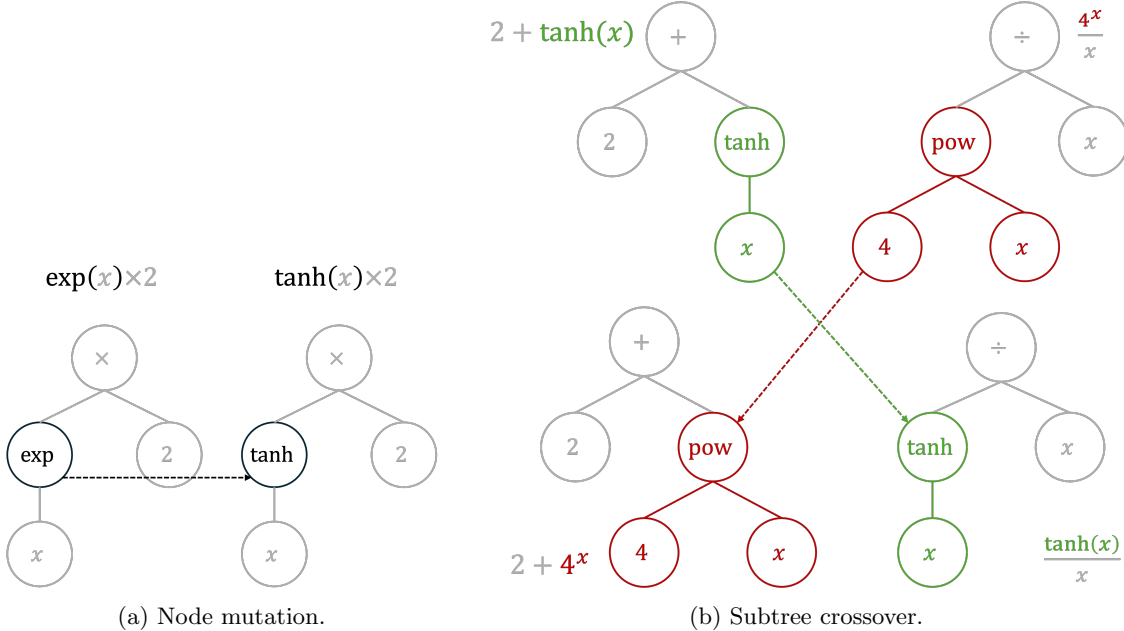


Figure 1: Genetic programming approach to symbolic regression. Functions are represented by expression trees. New functions are generated through mutation of tree nodes (left) and crossover between subtrees (right).

2. Challenges in traditional methods

2.1. Scenario 1: signal and background modeling

When analyzing proton-proton collision data at the LHC in the search for new physics signatures, the data are typically binned and presented as histograms representing physical observables, such as the invariant mass of the final-state particles. Each bin records the observed or expected number of collision events with mass values within that range. To search for new physics signatures, which are often hypothesized as narrow and small peaks over a smoothly falling background in the invariant mass distribution, parametric functions are required to model both the signal and background based on these binned distributions. These models are then used to perform hypothesis testing.

In the traditional approach to parametric modeling, one typically relies on manually guessing the appropriate family of functions that might describe the shape of the distribution. Although these distributions often represent physical observables, they are usually obtained after applying a series of selection cuts on various variables, which can introduce arbitrary shape effects into the final distributions being modeled. As a result, these distributions generally do not have a known underlying true function, making it impossible to derive a suitable functional form from first principles and leaving empirical constructions as the only option. In some cases, when a suitable function cannot be found to describe the distribution, one is forced to compromise by adjusting the analysis strategy, say splitting the main distribution into multiple sub-distributions and fitting them separately. This empirical approach has been the standard strategy in the HEP community, requiring significant manual effort to craft a candidate function and iteratively fine-tune it.

For example, a search for new physics in high-mass trijet events performed by the CMS experiment [8] modeled the background by fitting the trijet invariant mass distribution, m_{ijj} , using

three families of empirical functions. One of these functions takes the form:

$$f(x; N) = \frac{p_0(1-x)^{p_1}}{x^{\sum_{i=2}^N p_i \log^{i-2}(x)}}, \quad (1)$$

where $x = m_{\text{jj}}/\sqrt{s}$ is a dimensionless variable, p_i are free parameters, and N is a hyperparameter for the function form. The function was fitted multiple times with different trial N values, and the optimal value was determined through a separate statistical test, such as an F-test [9].

The challenge lies in the need to empirically craft a specific functional form, such as Eq. 1, for each individual distribution. These empirical functions are tailored to the particular distribution being fitted, making them rigid and potentially ineffective if there are slight changes in the data. For instance, variations in final-state objects, event selection strategies, or detector conditions during data collection can all introduce arbitrary shifts in the shape of the distribution. In such cases, function families that worked for past datasets may no longer be effective for future datasets, even within the same analysis channel, and an empirical searching for suitable functions must be repeated.

Analyses at the LHC have traditionally relied on this empirical fitting method when modeling signal and background processes from binned data. Examples include the milestone analyses that led to the discovery of the Higgs boson in 2012 [10–12], as well as some recent results from CMS searches for high-mass resonances in dijet [13], paired-dijet [14], trijet [8], diphoton [15], and dimuon [16] events.

In this context, SR has the potential to transform the approach to parametric modeling. By conducting a *machine-search* for suitable functional forms, SR significantly reduces the manual effort required in the modeling process, providing a more efficient and adaptive alternative to traditional methods.

2.2. Scenario 2: derivation of smooth descriptions from binned data

When predicting signal and background processes using simulation, there is always some degree of mismatch with the observed data, which may result from inaccuracies in theoretical predictions, mis-modeling of detector effects, or measurements errors. These discrepancies are corrected by applying data-to-simulation scale factors to the simulated events, ensuring that the simulation provides a more accurate representation of the observed data. Examples include jet energy scale corrections parameterized by jet p_T and η [17], heavy-flavor jet tagging efficiency corrections parameterized by the jet p_T and η [18], hadronic tau identification efficiency corrections parameterized by the tau p_T , η , and decay modes [19].

These scale factors are typically derived from binned data and applied as binned weights, resulting in coarse-grained corrections. When smooth scale factors are desired, the process often follows the same empirical approach, facing the same limitations discussed earlier. In cases where the scale factor is parameterized by more than one variable, it becomes even more challenging to empirically construct a suitable functional form, leading to a reliance on coarse-grained corrections.

Another common scenario involves data-driven background estimation methods, where transfer factors are derived to estimate events in the signal region based on those in the sideband region. For example, in a search for a boosted Higgs boson decaying to b quarks performed by the CMS experiment [20], the QCD multijet background was estimated from observed data by parameterizing the transfer factor, empirically constructed as the sum of products of Bernstein polynomials:

$$f(x_0, x_1) = \left(\sum_{\mu=0}^{n_{x_0}} \sum_{\nu=0}^{n_{x_1}} a_{\mu,\nu} b_{\mu,n_{x_0}}(x_0) b_{\nu,n_{x_1}}(x_1) \right) \times g(x_0, x_1), \quad (2)$$

‡Common coordinate system used to define particle kinematics in collider physics: p_T is transverse momentum and η is the pseudorapidity angle.

where $b_{n,N}(x)$ is the n -th Bernstein basis polynomial of degree N , $a_{\mu,\nu}$ are parameters to be extracted from a fit to observed data, and $g(x_0, x_1)$ is a function fitted separately to simulated events. The degrees of the Bernstein polynomials, n_{x_0} and n_{x_1} , are determined using an F-test.

By using SR, these empirical steps for deriving smooth scale factors can be significantly simplified into a single SR fit, eliminating the need to predefine a candidate function.

2.3. An alternative method: Gaussian process regression

An alternative fitting method is Gaussian process regression (GPR), which has been explored for these scenarios [21–23]. GPR models the dependent variable as following a Gaussian distribution at each point along the independent variable. The smoothness of the probability function is controlled by a chosen covariance kernel between bins. As a result, GPR provides a probabilistic prediction, yielding both a smooth mean function and a variance function, defining a very generic distribution of functions that describe the data, rather than a single exact function.

Fitting a GPR model to n data points requires inverting an $n \times n$ covariance matrix, which scales with a time complexity of $\mathcal{O}(n^3)$ [24]. This can become computationally prohibitive, especially for datasets with more than one independent variable. Additionally, integrating Bayesian GPR outputs into the standard HEP search framework requires subtle treatments [21], whereas SR directly provides explicit function templates that can be straightforwardly integrated into existing workflows.

Despite the potential of alternative methods like GPR, the HEP community continues to rely primarily on empirical approaches to parametric modeling. There is currently a lack of an efficient framework or package based on an alternative method that can be readily used out-of-the-box.

3. Proposed solution with symbolic regression

For the scenarios discussed above, we propose using SR to replace traditional methods, shifting the paradigm of parametric modeling in HEP.

In this paper, we introduce a Python API§ for the `SymbolFit` framework, which interfaces with `PySR` [7] (a high-performance SR library) and `LMFIT` [25] (a nonlinear least-square minimization library), aimed at *automating* parametric modeling of binned data using SR. We demonstrate the effectiveness of the framework in two common HEP applications: parametric modeling of signal and background, and the derivation of smooth scale factors. These applications are validated using five real datasets from new physics searches at the CERN LHC, along with several toy datasets. The key features of the framework are summarized below.

- **Pre-determined functional forms are no longer needed.** With SR, only minimal *constraints* are required to define the function space, such as specifying the allowed mathematical operators (+, ×, /, pow, sin(·), exp(·), etc.). This does not demand extensive and prior knowledge of the specific functions that describe particular distributions. The search for suitable functions is automatically performed by machine, eliminating the empirical and manual process typical of traditional methods. We show that a simple SR configuration can effectively fit a wide variety of distribution shapes.
- **Generating multiple candidate functions per fit.** SR based on genetic programming generates and evolves successive generations of functions, producing a batch of candidate functions in each search iteration. The same search configuration can be repeated with different random seeds to explore different suitable functions from the vast function space. This flexibility in generating a variety of candidate functions allows for adaptability across different downstream tasks.
- **Including uncertainty measure.** While SR algorithms alone are dedicated to function searching and do not inherently provide uncertainty estimates, our framework bridges the gap

§<https://github.com/hftsoi/symbolfit>

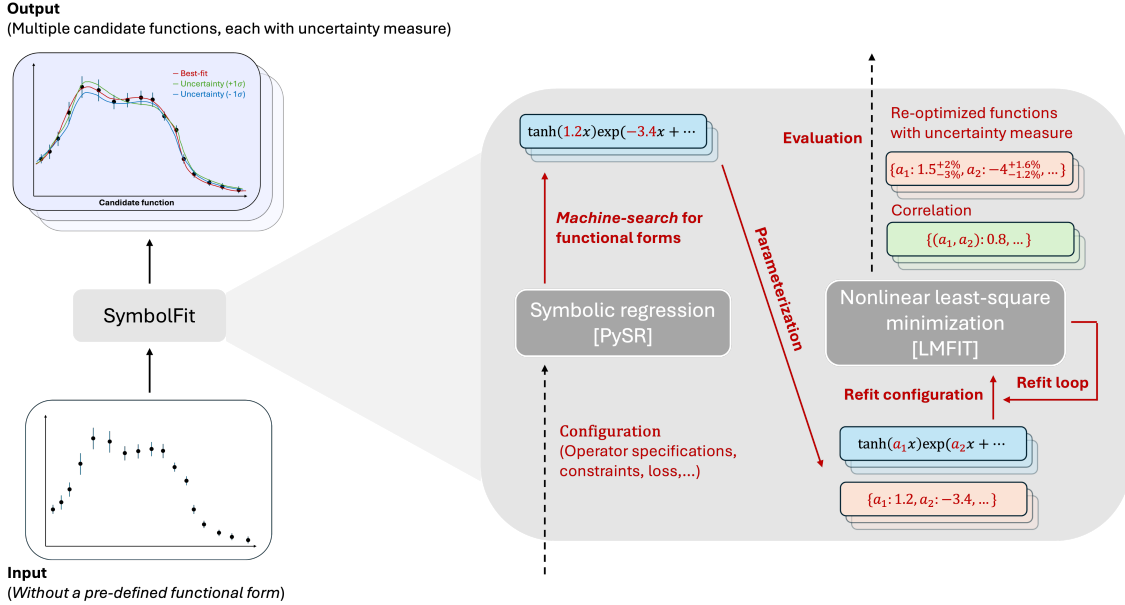


Figure 2: A schematic sketch of the internal steps within the SymbolFit framework illustrates how it interfaces with PySR and LMFIT to automate parametric modeling using SR. The process begins with an input dataset that does not require a predefined functional form. Functional forms are generated using SR, parameterized, and then re-optimized through standard nonlinear least-square minimization. The output is a batch of candidate functions, each with associated uncertainty estimates.

by incorporating a re-optimization process for the candidate functions. This step improves the best-fit models and generates uncertainty estimation.

- **Modeling of multi-dimensional data.** The framework easily accommodates modeling data with multiple variables, which is particularly useful in HEP scenarios where scale factors are sometimes parameterized by more than one variable.

Moreover, the framework is designed to automate the process as much as possible, minimizing manual effort. Results are automatically saved in easily readable formats, including visualizations for better interpretation. The candidate functions generated can be seamlessly integrated into downstream statistical inference tools commonly used in HEP, such as **Combine** [26] and **pyhf** [27, 28], as the output objects are in identical representation to those from traditional methods—closed-form functions. The efficiency of SR in generating a wide range of well-fitted functions per fit also allows flexible modeling as the function choice can be treated as a source of systematic uncertainty through the discrete profiling method [29].

The rest of the paper is structured as follows. Sec. 4 describes the SymbolFit framework. Sec. 5 presents demonstrations using a real LHC dataset as well as several toy datasets. Sec. 6 provides a summary of the work. More demonstrations are presented in **Appendix A**.

4. Method

The SymbolFit framework is illustrated in Fig. 2 and explained in the following.

- Input data.** We consider the input dataset $\{(\mathbf{x}^i, y^i, y_{\text{up}}^i, y_{\text{down}}^i)\}_{i=1}^n$, where \mathbf{x}^i represents one or more independent variables, y^i is the dependent variable with associated uncertainties $y_{\text{up/down}}^i$ at one standard deviation, and n is the number of data points.

In the context of binned histograms, which are commonly used in HEP data analysis, there are n bins. Here, \mathbf{x}^i represents the center of the i -th bin, and y^i is the bin content, representing the number of events within the bin. The associated uncertainties $y_{\text{up/down}}^i$ account for measurement errors or modeling inaccuracies.

- (ii) **Symbolic regression.** The core of the framework is to leverage SR to perform a *machine-search* for suitable functions to model the data, eliminating the need to predefine a functional form before fitting. We utilize PySR [7], a Python library for genetic programming-based SR, which is highly configurable in defining the function space for the search. The configuration process is highly simplified, requiring only the specification of allowed mathematical operators (+, ×, /, pow, sin(·), exp(·), etc.) and the constraints for the functional form. The objective of the search is to minimize:

$$\chi^2 \equiv \sum_{i=1}^n \left(\frac{f(\mathbf{x}^i) - y^i}{y_{\text{up}}^i \mathbf{1}_{f(\mathbf{x}^i) - y^i \geq 0} + y_{\text{down}}^i \mathbf{1}_{f(\mathbf{x}^i) - y^i < 0}} \right)^2, \quad (3)$$

where f is the candidate function. Since PySR uses a multi-population strategy to evolve and select functions, each run generates a batch of candidate functions. These functions are then re-optimized in subsequent steps to improve the fit and provide uncertainty estimates.

- (iii) **Parameterization.** SR algorithms search for exact functions but do not inherently provide uncertainty measures. However, uncertainty estimation is essential in HEP data analysis to gauge the reliability of the observation and prediction. To address this, we freeze the functional forms found by SR and then re-optimize all constants in each function using standard nonlinear minimization techniques. The uncertainties in these re-optimized constants are used as the uncertainty measure for the candidate functions.

First, within each candidate function, the constants are automatically identified and parameterized as $\{a_1, a_2, \dots\}$, with the original values stored as initial values for the re-optimization process.

- (iv) **Re-optimization fit (ROF).** To perform ROF of the functions, we utilize LMFIT, a nonlinear least-square minimization library, to perform a second-fit for the parameters while keeping the functional forms fixed. The objective is to minimize χ^2 defined in Eq. 3. The parameterized functions are parsed to identify the set of parameters to be varied, and initially, all parameters are allowed to vary in the fit.

In some cases, the minimization may fail to converge due to a too complex objective function. To handle these cases, a loop for ROF is implemented in the framework. This loop progressively reduces the number of degrees of freedom (NDF) by freezing more parameters to their initial values until the fit succeeds and all relative errors are below a pre-defined threshold.

Finally, the candidate functions are evaluated and ranked in the outputs.

To summarize, `SymbolFit` automates all these steps in the modeling process, including the computation of various goodness-of-fit scores and the evaluation of correlations between the parameters. This integration streamlines the workflow and minimizes manual intervention while providing full information for downstream statistical analysis. The computation time of the workflow is primarily due to the search for functional forms, for which we utilize a highly optimized SR algorithm PySR. As a result, the process is not computationally intensive and can be flexibly configured.

5. Demonstrations

We demonstrate the effectiveness of our framework using five real LHC datasets from new physics searches published recently, as well as several toy datasets.

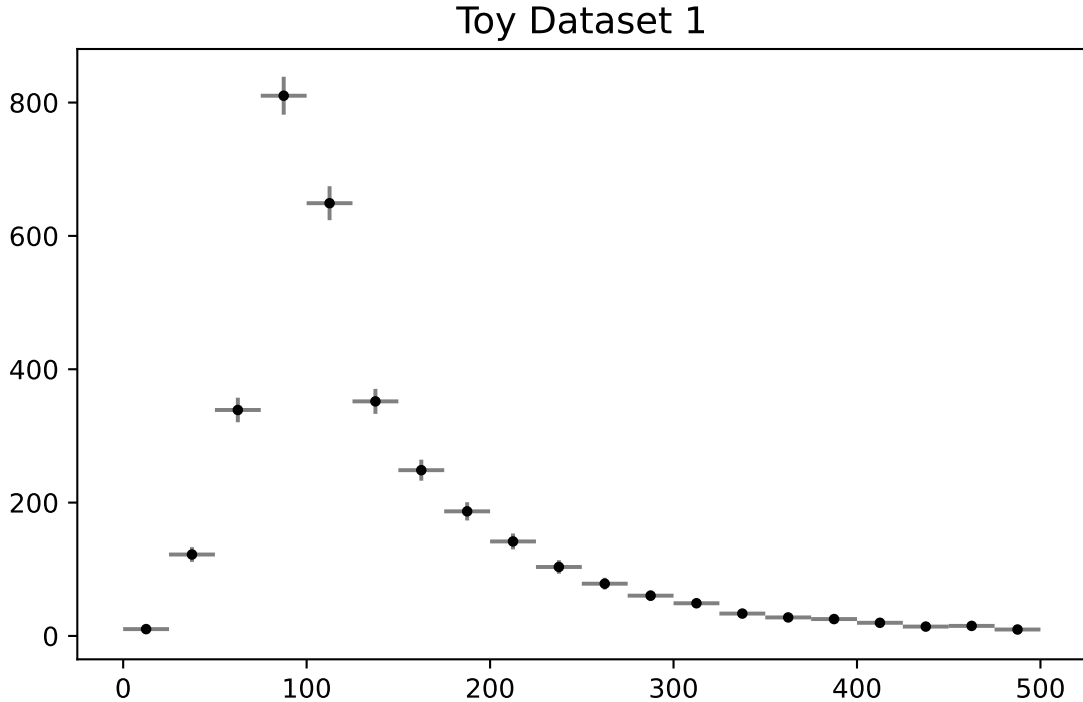


Figure 3: Toy Dataset 1: a 1D binned dataset with uncertainties represented by vertical error bars. The data points are manually generated without reference to an underlying function.

The LHC datasets consist of real proton-proton collision data at a center-of-mass energy of $\sqrt{s} = 13$ TeV, collected by the CMS experiment during Run 2. These datasets cover various search channels: dijet [13] in Sec. 5.2, diphoton [15] in Appendix A.2, trijet [8] in Appendix A.3, paired-dijet [14] in Appendix A.4, and dimuon [16] in Appendix A.5. Each dataset consists of 1D binned data of the invariant mass of the respective objects, where smooth background predictions are obtained through parametric modeling and then tested for excess events indicative of new physics. In all these analyses, CMS reported no evidence of new physics is observed in the data. Therefore, for our demonstrations, we assume that each invariant mass spectrum contains no signal. We also perform experiments to validate the SR outputs for signal extraction in these LHC datasets, and details of these steps are given in Sec. 5.2.

In addition to the real LHC datasets, we also generate toy datasets for demonstration purposes. While SR has been shown to successfully identify the correct underlying function from noisy data [1, 7], here we focus on toy data generated by hand without an underlying function to illustrate SR’s capability in modeling arbitrary distribution shapes. Toy Dataset 1, presented in Sec. 5.1, is a 1D binned distribution featuring a sharp peak and a high-end tail. Toy Dataset 2, presented in Sec. Appendix A.1, consists of three 1D binned distributions with various shapes. Toy Dataset 3, presented in Sec. 5.3, consists of three 2D binned distributions.

5.1. Toy Dataset 1 (1D) [signal modeling]

Fig. 3 shows Toy Dataset 1, which consists of a 1D binned distribution. The distribution has a shape characterized by a sharp peak and a high-end tail, and it is generated by hand without reference to an underlying function.

In spite of the lack of a true functional form, we demonstrate that our framework using SR

can replace the empirical process with minimal efforts. This can be seen from the simple Python snippet shown in List. 1¶, which defines the function space and configures PySR to perform a *machine-search* for functions to fit this dataset. In this example, the maximum function complexity is set to 60 to constrain the model size, ensuring that the total number of nodes in an expression tree does not exceed 60, with each operator equally weighted with a complexity of one. The allowed operators include two binary operators (+ and ×) and three unary operators (exp(·), tanh(·), and a custom-defined gauss(·) ≡ exp(−(·)²)). Constraints on operator nesting are imposed to prohibit scenarios like tanh(tanh(·)). The loss function used is χ^2 .

```

1 from pysr import PySRRegressor
2 import sympy
3
4 pysr_config = PySRRegressor(
5     model_selection = "accuracy",
6     timeout_in_seconds = 60*100,
7     niterations = 200,
8     maxsize = 60,
9     binary_operators = ["+", "*"],
10    unary_operators = [
11        "exp",
12        "gauss(x) = exp(-x*x)",
13        "tanh"
14    ],
15    nested_constraints = {
16        "tanh": {"tanh": 0, "exp": 0, "gauss": 0, "*": 2},
17        "exp": {"tanh": 0, "exp": 0, "gauss": 0, "*": 2},
18        "gauss": {"tanh": 0, "exp": 0, "gauss": 0, "*": 2},
19        "*": {"tanh": 1, "exp": 1, "gauss": 1, "*": 2}
20    },
21    extra_sympy_mappings={"gauss": lambda x: sympy.exp(-x*x)},
22    loss = "loss(y, y_pred, weights) = (y - y_pred)^2 * weights"
23 )

```

Listing 1: The Python code snippet that configures PySR to search for candidate functions for Toy Datasets 1 and 3.

We run `SymbolFit` with the PySR configuration shown in List. 1. This single fit generates 46 candidate functions (labeled #0 to #45), with function complexity values ranging from 1 to 60. These complexity values provide a rough estimate of the model size and are computed before any algebraic simplification. Four goodness-of-fit scores, including the χ^2/NDF and p-value, are plotted against function complexity in Fig. 4. As expected, more complex functions tend to offer better fits to the data, as their higher complexity makes them more expressive. The variety of functions produced provides flexibility, allowing one to choose a candidate function based on the desired fit quality for downstream tasks. For example, nine of the 46 candidate functions are listed in Tab. 1, showing χ^2/NDF values above 1, around 1, and below 1, offering a range of fit qualities

In general, the χ^2/NDF score improves significantly after ROF compared to the original function found from SR. This is because SR algorithms are typically focused on finding optimal functional forms rather than fine-tuning a specific function. As a result, the ROF step improves the constants in each function to achieve a better fit and provides uncertainty estimates for the parameters.

As an example, Fig. 5 shows candidate function #27. This candidate function has six parameters, with uncertainty variations for each plotted separately. The correlation matrix for these parameters is shown in Fig. 6.

To compare various candidate functions obtained from a single fit in the framework, Fig. 7 shows four candidate functions with a range of fit scores from very low to very high. This demonstrates that within the variety of functions generated per fit, there is a convergence from poorer to better-fitted functions, providing flexibility in choices. To illustrate the uncertainty

¶The option definitions can be found at <https://github.com/MilesCranmer/PySR> and Ref. [7].

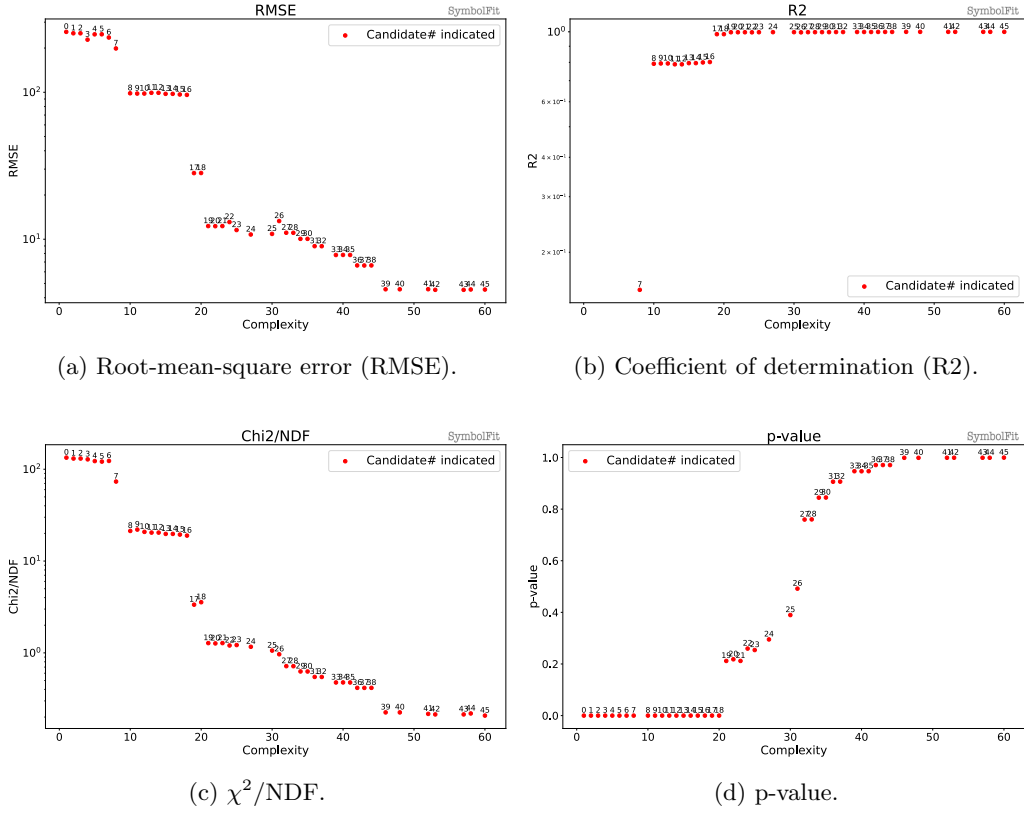


Figure 4: Goodness-of-fit scores vs. function complexity. A total of 46 candidate functions (labeled #0–#45) were obtained from a single fit on Toy Dataset 1 in Sec. 5.1. Candidate functions #10, #17, #27, and #38 are listed in Tab. 1, and their combined uncertainty coverage is presented in Fig. 7. Specifically, individual parameter variations for candidate function #38 are shown in Fig. 5, with its parameter correlation matrix shown in Fig. 6.

coverage for a candidate function, an ensemble of functions is generated by sampling parameters from a multidimensional normal distribution, considering their best-fit values and covariance matrix. The 68% quantile range of the function ensemble is shown for each candidate function, representing the total uncertainty coverage by simultaneously accounting for uncertainties in all parameters. These computations are all automatically performed within the SymbolFit framework.

Table 1: Nine examples of the 46 candidate functions obtained from a single fit for Toy Dataset 1. These functions were fitted to a scaled dataset (to enhance fit stability and prevent numerical overflow), which can be rescaled to describe the original dataset using the transformation: $f(x) \rightarrow 165 \times f(0.00211(x - 12.5))$. The comparison between the χ^2/NDF scores before and after the ROF step is presented. The total uncertainty coverage of candidate functions #10, #17, #27, #38 is shown in Fig. 7. Individual parameter variations for candidate function #27 are plotted in Fig. 5. The function complexity values, providing a rough estimate of the model size, are computed before algebraic simplification. Numerical values are rounded to three significant figures for display purposes.

Complexity	Candidate function (after ROF)	# param.	χ^2/NDF (before ROF)	χ^2/NDF (after ROF)	p-value (after ROF)
12 (#10)	$0.101 + 23.3\text{gauss}(-3.74x)x$	2	374.3 / 18	374.3 / 18	10^{-68}
19 (#17)	$0.061 + (5.11x + 5.11\text{gauss}(-2.34 + 12.6x))\text{gauss}(2.52x)$	4	54.06 / 16	53.59 / 16	10^{-06}
22 (#20)	$0.0837 + (4.76\text{gauss}(-15.7x + 2.84) + 10.4 \tanh(x))\text{gauss}(2.99x)$	6	48.47 / 14	17.76 / 14	0.218
27 (#24)	$(4.79\text{gauss}((-5.5 + 2x)(-0.538 + 3x)) + 10.2x)\text{gauss}(-3.03x) + 0.0841$	6	16.88 / 14	16.31 / 14	0.2951
31 (#26)	$(4.9x + 4.9\text{gauss}(-2.79 + 15.4x) + 4.9 \tanh(x))\text{gauss}(3x) + 0.0789\text{gauss}(x) \exp(x)$	4	15.79 / 16	15.45 / 16	0.4919
32 (#27)	$(5.13\text{gauss}(-16.7x + 3.05) + 13.1x)\text{gauss}(x(-4.68 + x) + x) + 0.0661$	6	12 / 14	10.04 / 14	0.760
37 (#32)	$(5.07\text{gauss}((-4.42 + 2x)(-0.724 + 4x)) + a5 \tanh(x) + 7.79x) \times \text{gauss}(x(-4.65 + x) + x) + 0.066$	6	8.359 / 14	7.655 / 14	0.9065
44 (#38)	$(5.08\text{gauss}((-4.7 + 4x)(-0.719 + 4x)) + 12.7x) \times \text{gauss}(x(-4.66 + x) + x) + 0.0662$	6	6.278 / 14	5.826 / 14	0.971
52 (#41)	$0.0657 + (5\text{gauss}((-4.96 + 6x)(-0.712 + 4x)) + 12.4 \tanh(x)) \times \text{gauss}(x(-4.6 + x) + x) - 0.00624x$	6	3.564 / 14	3.032 / 14	0.999

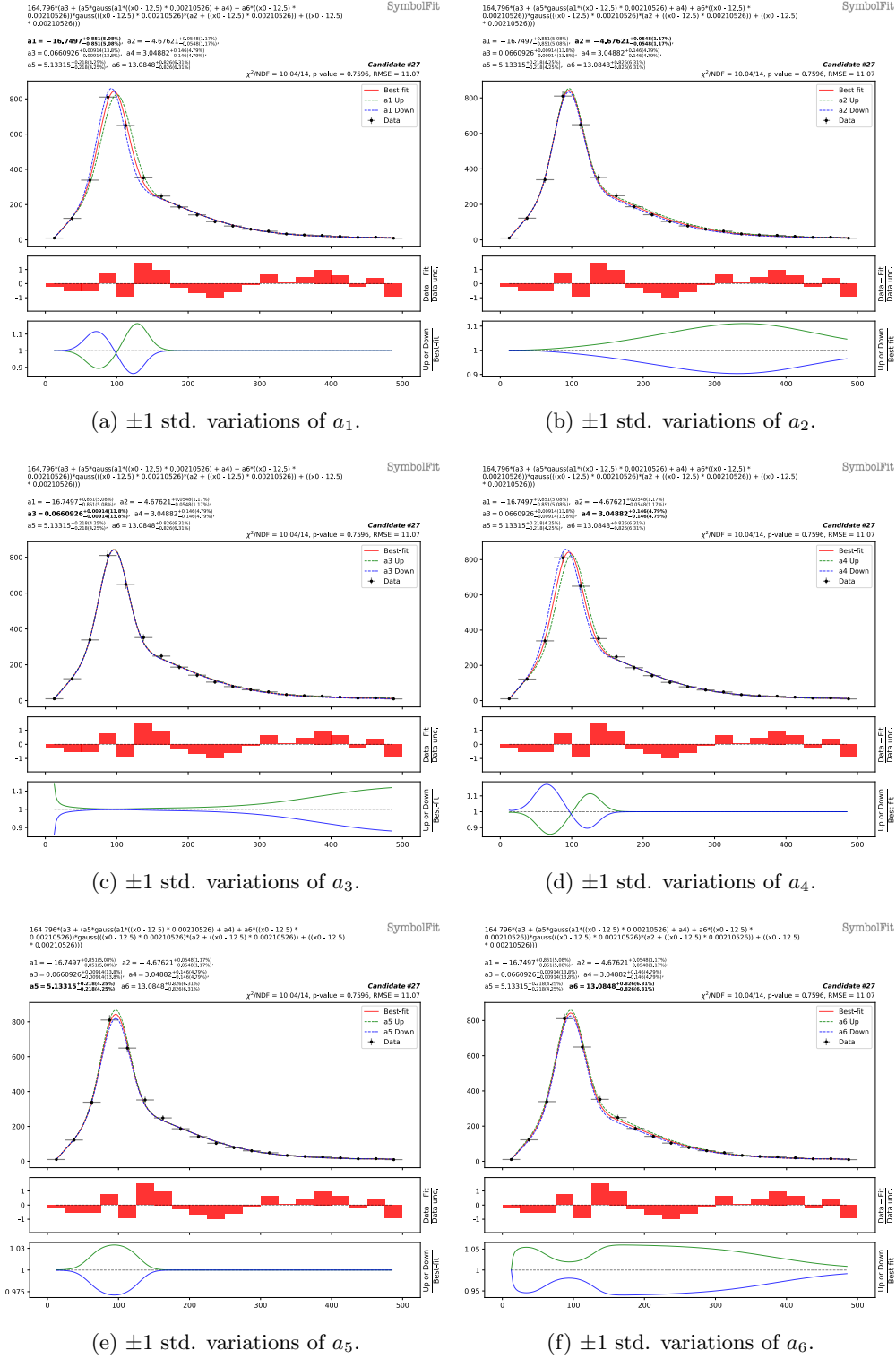


Figure 5: Individual parameter variations in candidate function #27 from a fit to Toy Dataset 1. The parameterized form of this function is shown at the top of each subfigure, along with the best-fit values of the parameters and their associated uncertainties. Each subfigure shows the same function, but with one parameter shifted by its ± 1 standard deviation (green/blue), while the other parameters remain fixed at their best-fit values. The function with all parameters held at their best-fit values is plotted in red and compared to the data, represented by black data points. The middle panel shows the weighted residual error: $\frac{\text{Data}-\text{Fit}}{\text{Data unc.}}$. The bottom panel shows the ratio of the function with the uncertainty variations to the best-fit function.

$$164.796*(a3 + (a5*gauss(a1*((x0 - 12.5) * 0.00210526) + a4) + a6*((x0 - 12.5) * 0.00210526))*gauss((x0 - 12.5) * 0.00210526)*(a2 + ((x0 - 12.5) * 0.00210526)) + ((x0 - 12.5) * 0.00210526))$$

SymbolFit

$$a1 = -16,7497^{+0.851(5.08\%)}_{-0.851(5.08\%)}, a2 = -4,67621^{+0.0548(1.17\%)}_{-0.0548(1.17\%)},$$

$$a3 = 0,0660926^{+0.00914(13.8\%)}_{-0.00914(13.8\%)}, a4 = 3,04882^{+0.146(4.79\%)}_{-0.146(4.79\%)},$$

$$a5 = 5,13315^{+0.218(4.25\%)}_{-0.218(4.25\%)}, a6 = 13,0848^{+0.826(6.31\%)}_{-0.826(6.31\%)}$$

Candidate #27

$$\chi^2/NDF = 10,04/14, p\text{-value} = 0,7596, RMSE = 11,07$$

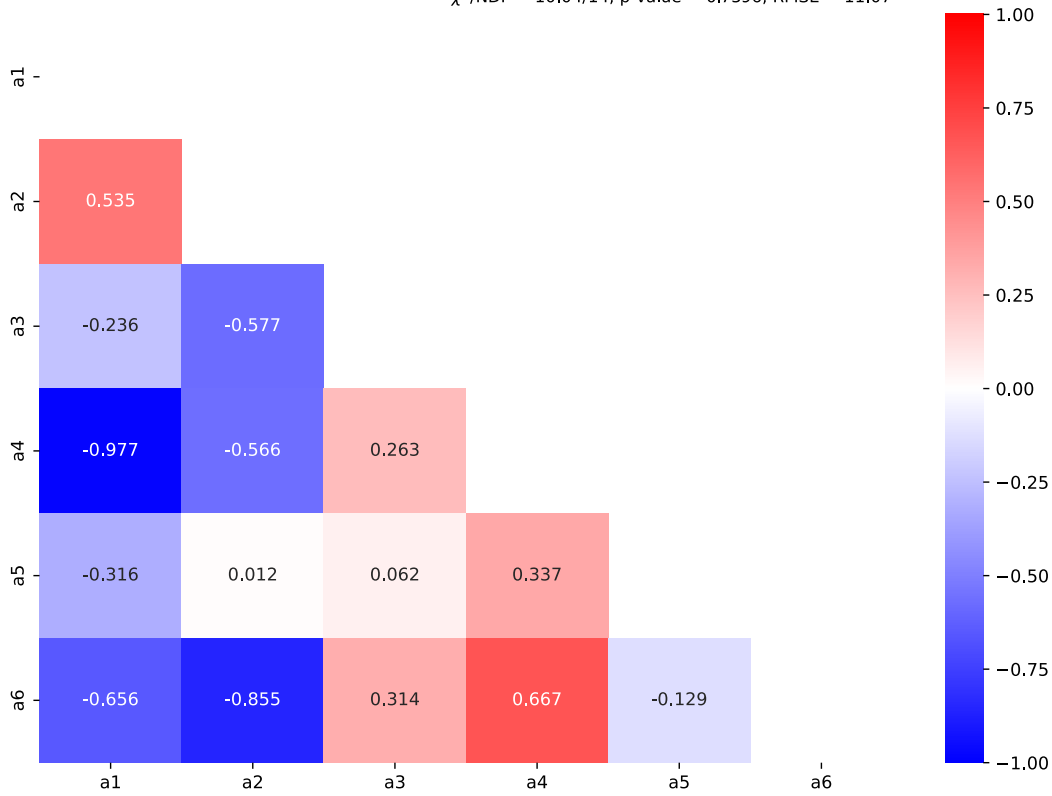
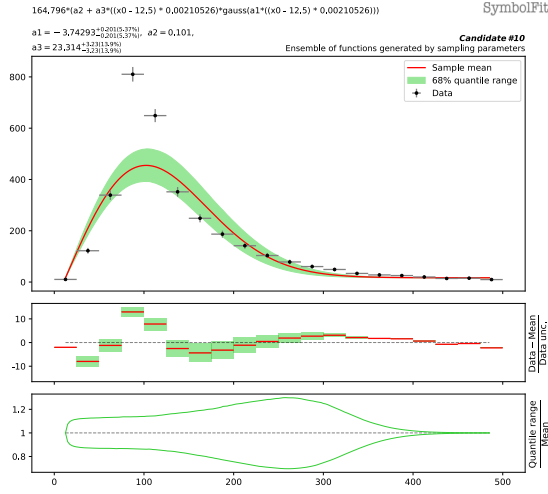
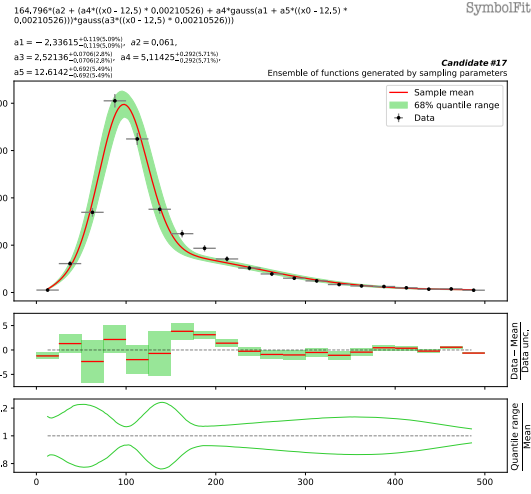


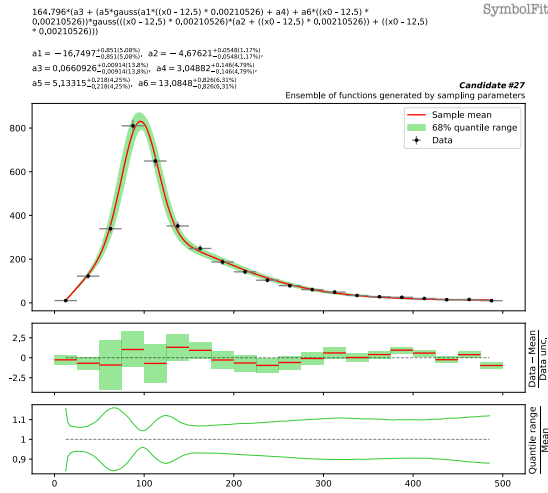
Figure 6: Correlation matrix for the parameters of candidate function #27 from a fit to Toy Dataset 1 (see Tab. 1 and Fig. 5).



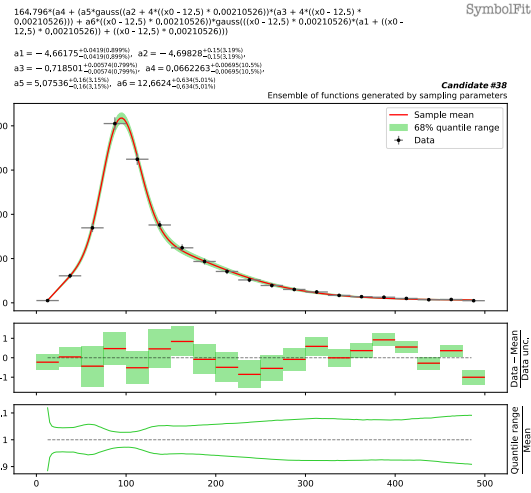
(a) Candidate function #10.



(b) Candidate function #17.



(c) Candidate function #27.



(d) Candidate function #38.

Figure 7: Convergence of candidate functions to the data (Toy Dataset 1), from lower to higher function complexity values. To visualize the total uncertainty coverage of each candidate function, the green band in each subfigure represents the 68% quantile range of functions obtained by sampling parameters, taking into account the best-fit values and the covariance matrix within a multidimensional normal distribution. The red line denotes the mean of the function ensemble. At the top of each subfigure, the candidate function and the fitted parameters are shown. The middle panel shows the weighted residual error: $\frac{\text{Data} - \text{Mean}}{\text{Data unc.}}$. The bottom panel shows the ratio of the 68% quantile range to the mean.

5.2. CMS dijet dataset (1D) [background modeling]

CMS performed a search for high-mass dijet resonances using proton-proton collision data at a center-of-mass energy of $\sqrt{s} = 13$ TeV and reported no significant deviations from the Standard Model prediction [13]. The dataset for the dijet spectrum is publicly available on HEPDATA at Ref. [30]. In the analysis, CMS used an empirical 4-parameter function to model the background contribution in the distribution of the dijet invariant mass, m_{jj} :

$$f(x) = \frac{p_0(1-x)^{p_1}}{x^{p_2+p_3 \ln(x)}}, \quad (4)$$

where $x = m_{jj}/\sqrt{s}$ is dimensionless and $p_{\{0,1,2,3\}}$ are free parameters. While this function fits the current dijet spectrum reasonably well, it may be too rigid to accommodate potential future changes in the dijet spectrum due to modifications in analysis strategies or detector performance.

For our experiments, we use the PySR configuration shown in List. 2 to fit the dijet dataset. The main difference between List. 2 and List. 1 used in Sec. 5.1 is that it does not explicitly include a Gaussian operator, as the mass spectrum assumes no peaks in the background. This same SR configuration is also applied to the other LHC datasets as well as Toy Dataset 2, generating a range of well-fitted functions for each case despite their very different distribution shapes, demonstrating the flexibility of the SR approach.

```

1 from pysr import PySRRegressor
2
3 pysr_config = PySRRegressor(
4     model_selection = "accuracy",
5     timeout_in_seconds = 60*100,
6     niterations = 200,
7     maxsize = 80,
8     binary_operators = [
9         "+", "*", "/", "^"
10    ],
11    unary_operators = [
12        "exp",
13        "tanh",
14    ],
15    nested_constraints = {
16        "exp": {"exp": 0, "tanh": 0, "*": 2, "/": 1, "^": 1},
17        "tanh": {"exp": 0, "tanh": 0, "*": 2, "/": 1, "^": 1},
18        "*": {"exp": 1, "tanh": 1, "*": 2, "/": 1, "^": 1},
19        "^": {"exp": 1, "tanh": 1, "*": 2, "/": 1, "^": 0},
20        "/": {"exp": 1, "tanh": 1, "*": 2, "/": 0, "^": 1},
21    },
22    loss="loss(y, y_pred, weights) = (y - y_pred)^2 * weights",
23 )

```

Listing 2: The Python code snippet configures PySR to search for candidate functions for the dijet dataset. This same configuration is used for the other four LHC datasets and Toy Dataset 2 with variations in the maximum complexity values.

To validate the SR approach in background modeling and signal extraction, we start with the original dijet spectrum and generate pseudodata by injecting a small and narrow Gaussian signal, $s_0 \frac{1}{\sqrt{2\pi}s_2} \exp\left(-\frac{(x-s_1)^2}{2s_2^2}\right)$, centered at $m_{jj} = 3.1$ TeV (s_1), with a width of 0.2 TeV ($2s_2$) and a signal strength of $s_0 = 10$. The injected signal intensity per bin is perturbed by 10% random noise.

To model the background, we first blind the signal region by masking the m_{jj} bins near the injected signal peak, specifically between 2.7 and 3.5 TeV, and then perform fits to this blinded pseudodata. We conduct three separate SymbolFit runs on the blinded pseudodata, using the same fit configuration but with different random seeds. This demonstrates that a variety of well-fitted functions can be obtained from the same fit configuration, given the vast function space in the SR search. From each of the three fits, one candidate function is selected, referred to as “SR

Table 2: The candidate functions are obtained from three fits using different random seeds, fitted to the pseudodata of the dijet spectrum with the (injected) signal region blinded. The fits were performed on a scaled dataset (to enhance fit stability and prevent numerical overflow), and the functions can be transformed back to describe the original spectrum using the transformation: $x \rightarrow 0.000145(x - 1568.5)$. These functions are plotted and compared with the blinded pseudodata in Fig. 9. Numerical values are rounded to three significant figures for display purposes.

	Candidate function (after ROF)	# param.	χ^2/NDF (before ROF)	χ^2/NDF (after ROF)	p-value (after ROF)
SR model 1	$(570x(x(-0.423 \exp(2x) + \exp(x)) + x) + 149) \times (0.00328 + 0.0304 \tanh(x))^{4.87x}$	5	400.5 / 30	29.21 / 30	0.507
SR model 2	$(145(0.958 + x)^{\tanh(-0.711+4.32x)} + 145 \tanh(x)) \times (0.00591 + 0.0522 \tanh(x))^{5.48x}$	5	103.8 / 30	29.91 / 30	0.47
SR model 3	$\text{pow}(149(0.0101x + 0.0101 \tanh(0.171 + 0.724x)), x + (2.38x \tanh(-0.71 + x) + 2.39) \tanh(x) + \tanh(x))$	5	214.8 / 30	30.93 / 30	0.419

model 1”, “SR model 2”, and “SR model 3”, respectively. These three background models are then compared with the empirical function in Eq. 4 used by CMS, referred to as the “empirical model (CMS)”.

Tab. 2 lists the three SR models, each obtained from a fit initialized with a different random seed. The χ^2/NDF scores improve significantly after the ROF step compared to the original functions returned by PySR, with the final scores close to 1. The three background models fit the blinded pseudodata well, as shown in Fig. 8 for the total uncertainty coverage and Fig. 9 for a comparison with the empirical model used by CMS.

Once the background models are established, we incorporate a parameterized Gaussian signal template into each model $f(x)$:

$$f(x) + s_0 \frac{1}{\sqrt{2\pi}s_2} \exp\left(-\frac{(x - s_1)^2}{2s_2^2}\right). \quad (5)$$

In the following analysis, when the model is fitted to the unblinded pseudodata with $s_0 = 0$ held fixed, it is referred to as a background-only (b-only) fit. When s_0 is allowed to vary, it is referred to as a signal-plus-background (s+b) fit.

Now, we unblind the pseudodata and perform b-only fits and s+b fits on the full dijet spectrum. Since the pseudodata contain an injected signal, we expect to observe an excess of events over the background model around the signal location, provided the background is properly modeled and not overly fitted. When performing the s+b fits, we expect that the excess of events observed in the b-only fits will diminish as the signal is accounted for by the model template. The results of the b-only and s+b fits for each model are compared and shown in Fig. 10. In all three SR models, as well as the CMS empirical model, the excess of events over the background around the injected signal location observed in the b-only fits is reduced in the s+b fits, demonstrating that the models are sensitive to the injected signal. Tab. 3 lists the χ^2/NDF scores for each model, indicating the fit performance in response to the presence of the injected signal.

Additionally, to assess whether the SR models can extract the injected signals, we generate multiple sets of pseudodata by injecting Gaussian signals with varying mean values between 2980 to 3150 GeV and signal strengths ranging from 2 to 38. We then perform the s+b fits to extract the corresponding signal parameters. Fig. A15 shows the extracted signal parameters (mass and strength) plotted against their injected values. All three SR models are capable of extracting the correct signal parameter values within reasonable uncertainties. They perform comparably to the

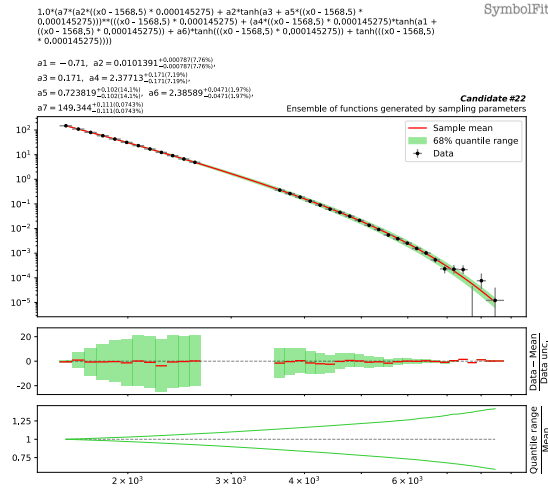
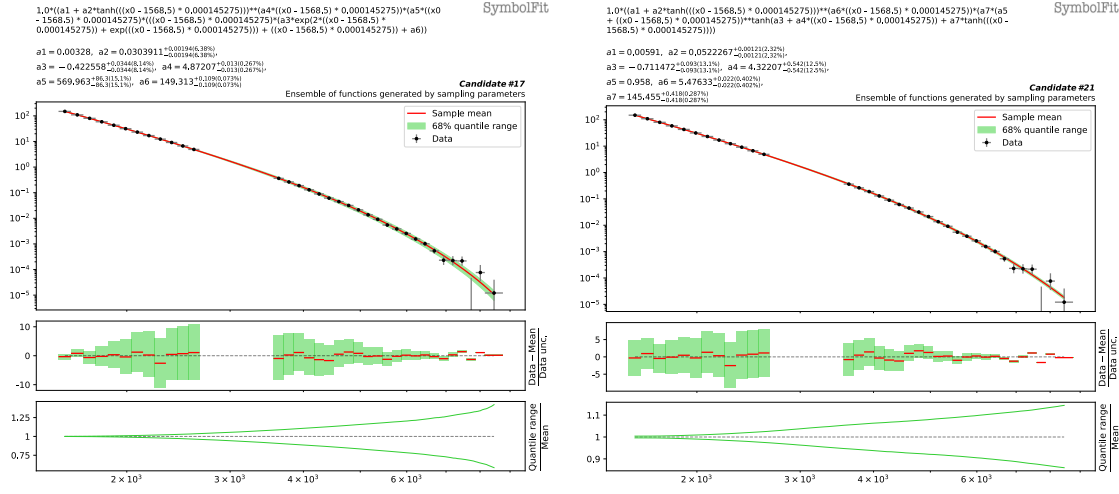


Figure 8: The three SR models fitted to the pseudodata of the dijet spectrum with the signal region blinded (see Tab. 2). To visualize the total uncertainty coverage of each candidate function, the green band in each subfigure represents the 68% quantile range of functions obtained by sampling parameters, taking into account the best-fit values and the covariance matrix within a multidimensional normal distribution. The red line denotes the mean of the function ensemble. At the top of each subfigure, the candidate function and the fitted parameters are shown. The middle panel shows the weighted residual error: $\frac{\text{Data} - \text{Mean}}{\text{Data unc.}}$. The bottom panel shows the ratio of the 68% quantile range to the mean.

empirical model used by CMS and, in some cases, yield more accurate fitted values, demonstrating that functions obtained from SR are effective for such tasks.

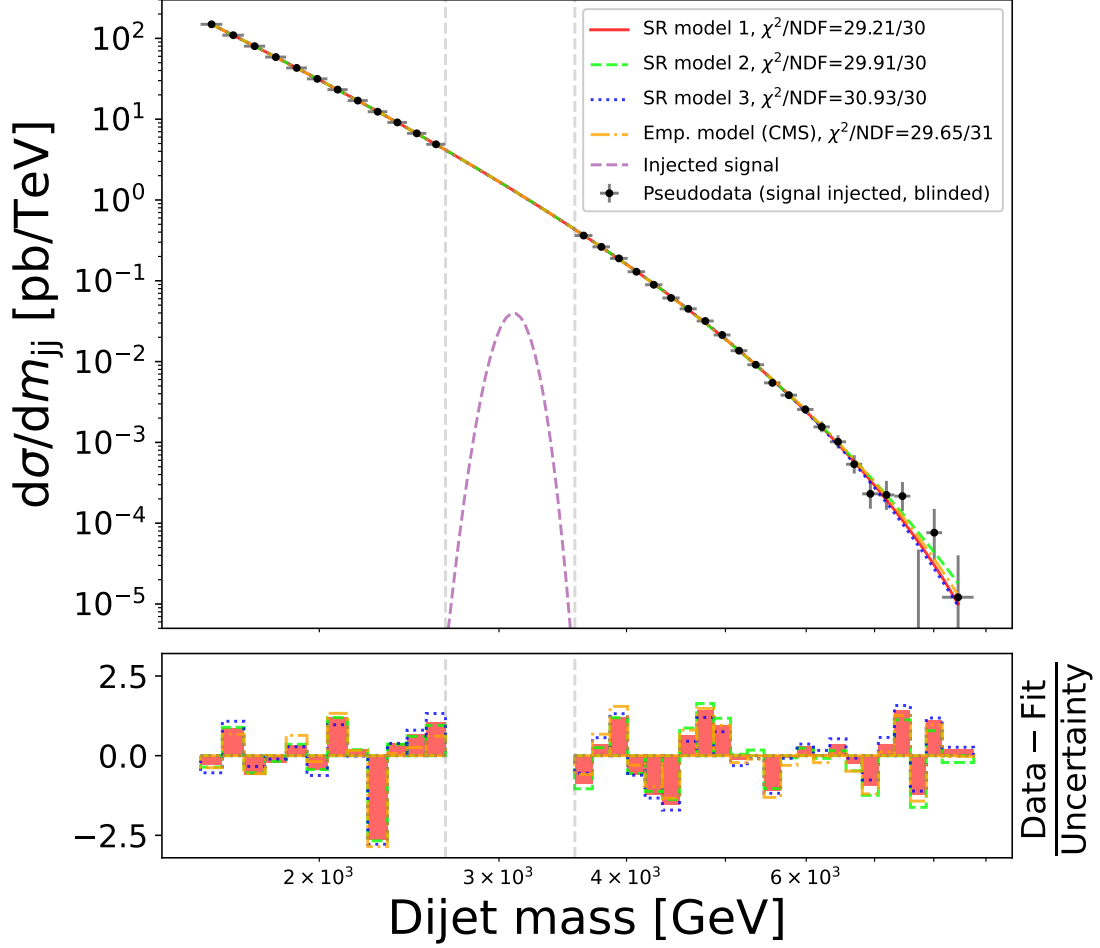


Figure 9: Pseudodata of the dijet spectrum with the injected signal shown in the blinded signal region. The three SR models (see Tab. 2) are compared against the empirical model used by CMS. The lower panel shows the residual error per bin, measured in units of the data uncertainty.

Table 3: Comparison of the χ^2/NDF scores from three types of fits to the dijet dataset: the b-only fits to the blinded pseudodata, the b-only fits to the unblinded pseudodata, and the s+b fits to the unblinded pseudodata. The background models used for the fits are listed in Tab. 2, and the fits are shown in Fig. 9 (blinded) and Fig. 10 (unblinded).

	χ^2/NDF (b-only, blinded)	χ^2/NDF (b-only, unblinded)	χ^2/NDF (s+b, unblinded)
SR model 1	29.21 / 30 = 0.974	54.48 / 37 = 1.47	36.19 / 35 = 1.03
SR model 2	29.91 / 30 = 0.997	51.25 / 37 = 1.39	36.28 / 35 = 1.04
SR model 3	30.93 / 30 = 1.03	65.33 / 37 = 1.77	39.98 / 35 = 1.14
Emp. model (CMS)	29.65 / 31 = 0.956	54.63 / 38 = 1.44	35.07 / 36 = 0.974

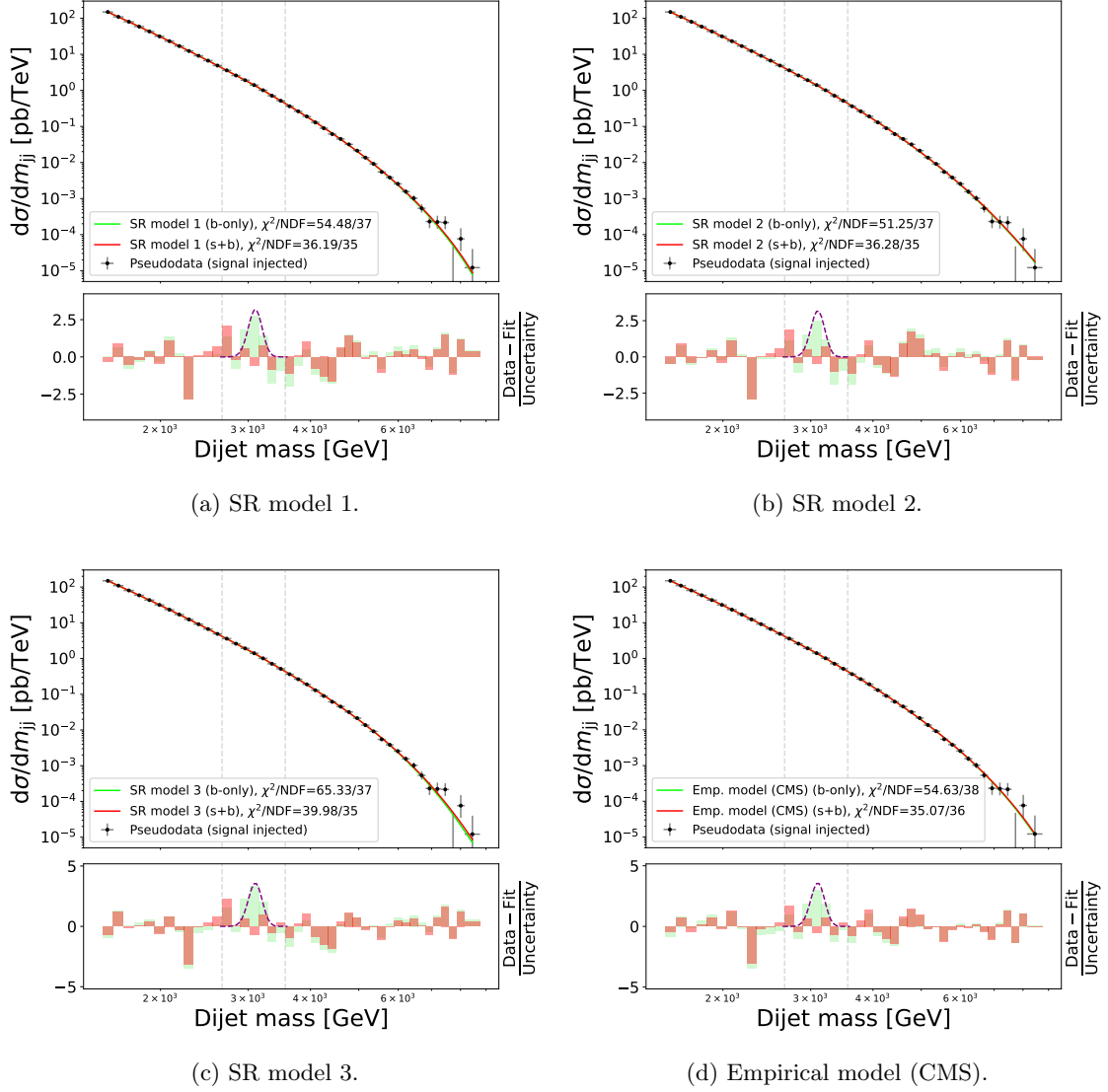
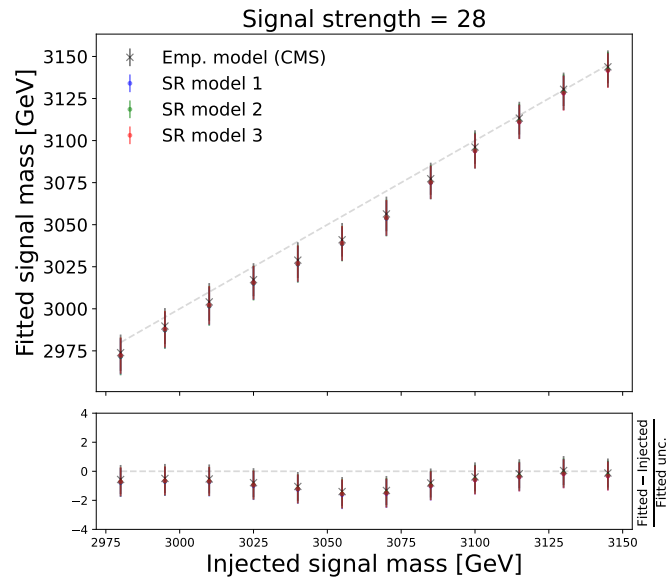
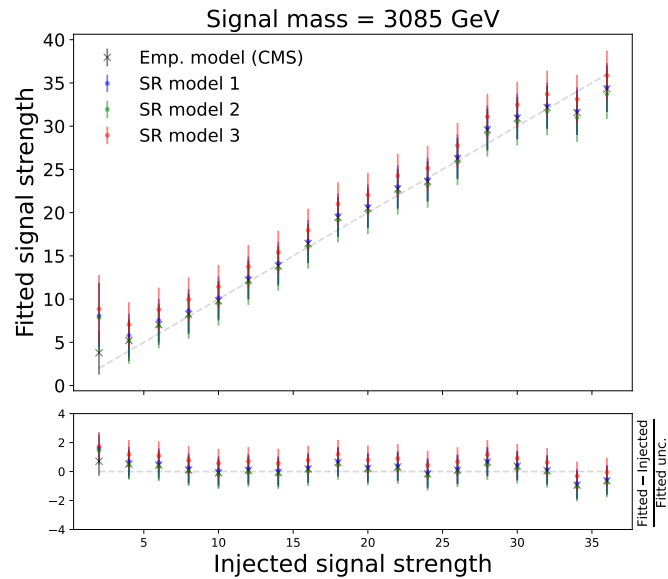


Figure 10: Comparison of the b-only fits and the s+b fits to the unblinded pseudodata of the dijet spectrum. The lower panel shows the residual error per bin, measured in units of the data uncertainty. The shape of the injected signal is also shown.



(a) Fitted vs. injected signal mass at a specified signal strength value.



(b) Fitted vs. injected signal strength at a specified signal mass value.

Figure 11: Fitted values vs. the true values of the parameters of the injected signal in the dijet dataset. The bottom panels show the residual error in units of the fitted uncertainty.

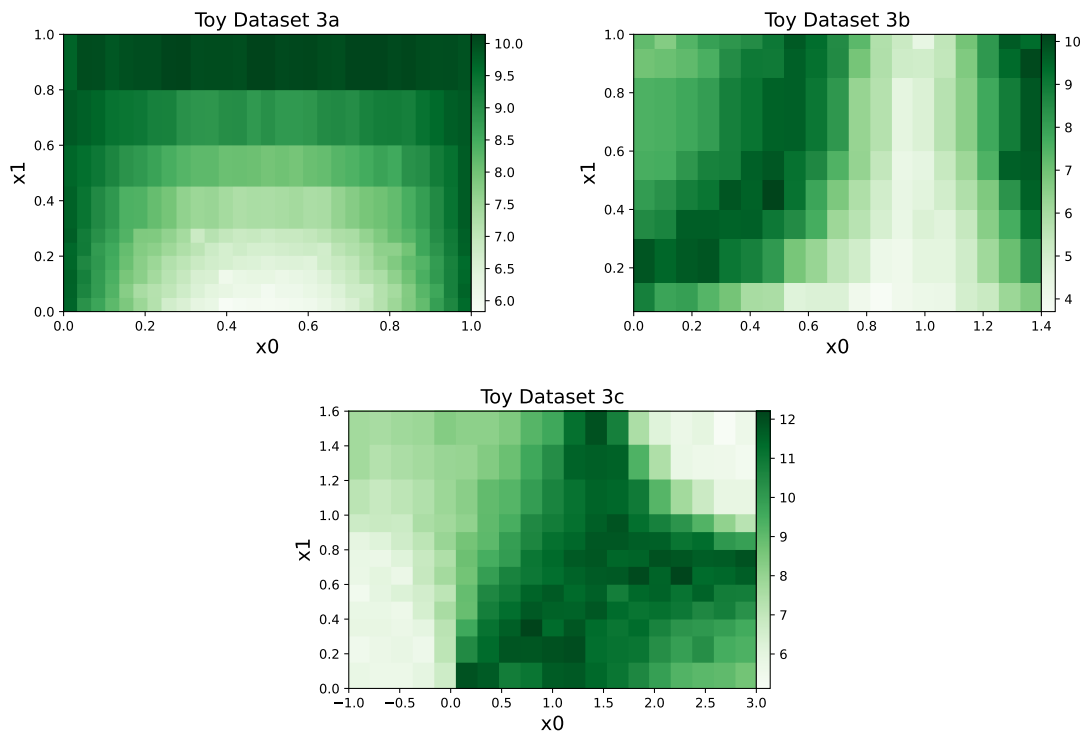


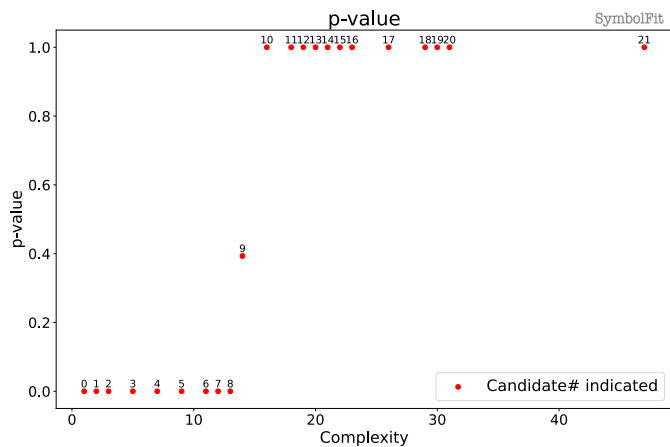
Figure 12: Toy Dataset 3: three 2D binned sub-datasets manually generated without reference to an underlying function.

5.3. Toy dataset 3 (2D) [arbitrary shapes]

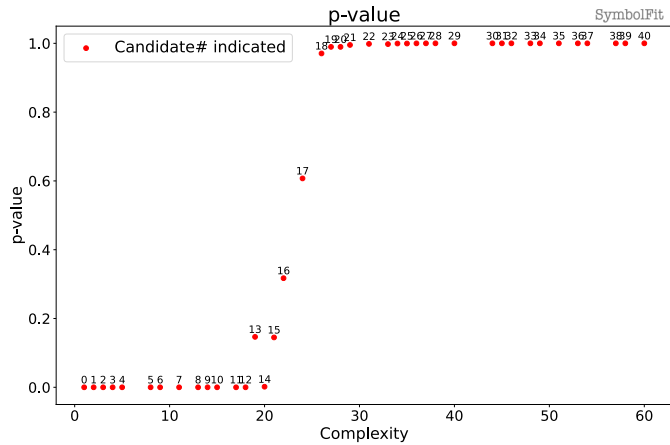
In Toy Dataset 3, we consider three 2D binned sub-datasets, labeled 3a, 3b, and 3c, as shown in Fig. 12. These datasets are manually generated without reference to an underlying function and are used to demonstrate applications such as deriving smooth scale factors from binned data with more than one independent variable. The framework can be easily extended to datasets with multiple independent variables.

We use the same PySR configuration applied to Toy Dataset 1, as shown in List. 1, to fit these 2D binned datasets. For each sub-dataset, a single run of SymbolFit is performed to generate a batch of candidate functions. Fig. 13 shows the p-value plotted against function complexity. Several candidate functions for each sub-dataset are selected and listed in Tab. 4.

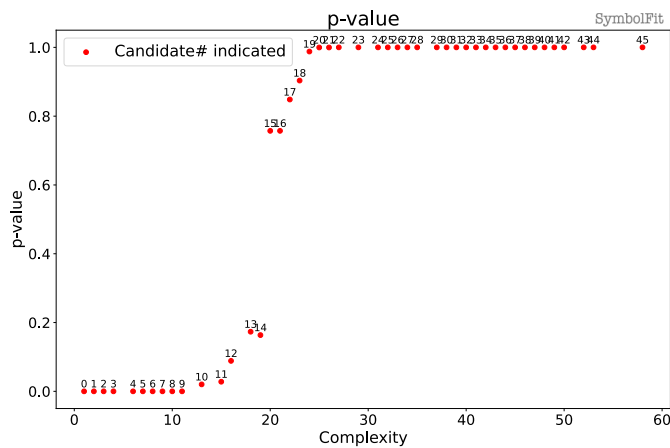
A candidate function is shown for each of the three sub-datasets: #26 for sub-dataset 3a in Fig. 14, #39 for sub-dataset 3b in Fig. 15, and #37 for sub-dataset 3c in Fig. 16.



(a) Toy Dataset 3a.



(b) Toy Dataset 3b.



(c) Toy Dataset 3c.

Figure 13: p-value vs. function complexity. A total of 22, 41, and 46 candidate functions (labeled #0–#21, #0–#40, and #0–#45) were obtained from a single fit on Toy Dataset 3a, 3b, and 3c, respectively.

Table 4: Examples candidate functions for Toy Dataset 3 are listed. The example candidate functions—#12 for Toy Dataset 3a, #38 for 3b, and #34 for 3c—are plotted in Fig. 14, Fig. 15, and Fig. 16, respectively. Numerical values are rounded to three significant figures for display purposes.

Toy Dataset 3a					
Complexity	Candidate function (after ROF)	# param.	χ^2/NDF (before ROF)	χ^2/NDF (after ROF)	p-value (after ROF)
14 (#9)	$8.08\text{gauss}(x_0(-2.76 + x_1)) + 9.46x_0 + x_1$	3	292.9 / 287 = 1.02	292.9 / 287 = 1.02	0.393
19 (#12)	$9.58 + 14.4x_0(-0.994 + x_0)(\text{gauss}(x_1))^2 + x_1$	3	41.06 / 287 = 0.1431	32.75 / 287 = 0.1141	1.0
20 (#13)	$9.92 + 16.0x_0(-1 + x_0)(-0.541x_1 + \text{gauss}(x_1))$	4	15.8 / 286 = 0.0552	13.22 / 286 = 0.0462	1.0
Toy Dataset 3b					
Complexity	Candidate function (after ROF)	# param.	χ^2/NDF (before ROF)	χ^2/NDF (after ROF)	p-value (after ROF)
22 (#16)	$1.38 \exp(x_0^2) + 7.08\text{gauss}(-0.736x_1 + x_0^2 + x_0) + \tanh(x_1)$	3	156.7 / 149 = 1.052	156.7 / 149 = 1.052	0.3172
24 (#17)	$1.3 \exp(x_0^2) + 6.83\text{gauss}(-0.726x_1 + x_0^2 + x_0) + \tanh(2.48x_1)$	3	144.4 / 149 = 0.969	143.7 / 149 = 0.964	0.607
26 (#18)	$0.345 + 5.73x_1\text{gauss}(x_1) + 5.73\text{gauss}(-0.809x_1 + 2.62x_0^2) + \exp(x_0^2)$	3	118.4 / 149 = 0.7945	118.1 / 149 = 0.7928	0.971
57 (#38)	$4.83\text{gauss}(0.571 + 1.18x_0(2x_0 + x_1) + x_1(-3.07 + x_1)) + 2\text{gauss}(x_1) - \tanh(-4.87 + 4.67x_0 + 2.74x_1) + 0.129 \exp(x_0(1.54 + x_0)) + 4.83 \tanh(x_1) + 0.335x_1$	8	52.8 / 144 = 0.3667	52.06 / 144 = 0.3615	1.0
Toy Dataset 3c					
Complexity	Candidate function (after ROF)	# param.	χ^2/NDF (before ROF)	χ^2/NDF (after ROF)	p-value (after ROF)
20 (#15)	$-0.886x_0(\text{gauss}(x_0x_1) + \exp(x_1)) + 6.74 \tanh(x_0) + 8.24$	3	214.2 / 225 = 0.9521	209.9 / 225 = 0.9328	0.7573
23 (#18)	$1.9x_0(-0.296 + \text{gauss}(x_1^2)) + 4.38\text{gauss}(-1.14 + x_0) + 6.56$	5	198.1 / 223 = 0.8881	196 / 223 = 0.8789	0.9035
25 (#20)	$4.33\text{gauss}(-1.14 + x_0) + 6.57 + x_0(-0.238 + \text{gauss}(x_1^2)) \times (1.33 + x_1)$	5	160.3 / 223 = 0.7188	157.3 / 223 = 0.7054	0.9997
42 (#34)	$5.16 + x_1 + (-0.229 + \text{gauss}(-0.481 + x_1^2(-1.25 + x_0))) \times (1.48x_0 + 1.48 \tanh(4.66x_0) + 2.35\text{gauss}(x_1)) + 3.53\text{gauss}(-0.907 + x_0)$	8	59.38 / 220 = 0.2699	54.89 / 220 = 0.2495	1.0

SymbolFit

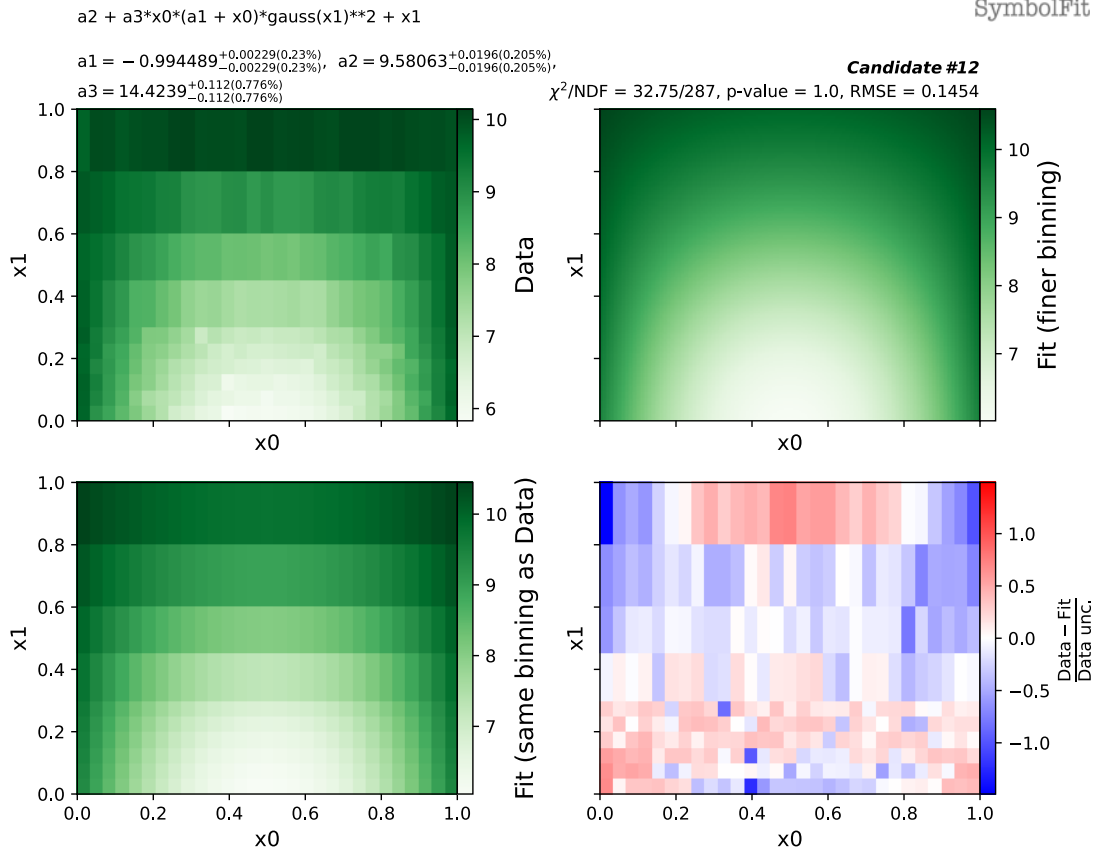


Figure 14: Candidate function #12 for Toy Dataset 3a (see Tab. 4). The parameterized form of this function is shown at the top of the figure, along with the best-fit values and associated uncertainties. Upper left: the binned data being fitted. Lower left: the candidate function plotted with the same binning as the fitting data. Upper right: the candidate function plotted with a finer binning. Lower right: the residual error, $\frac{\text{Data} - \text{Fit}}{\text{Data unc.}}$, in units of the data uncertainty.

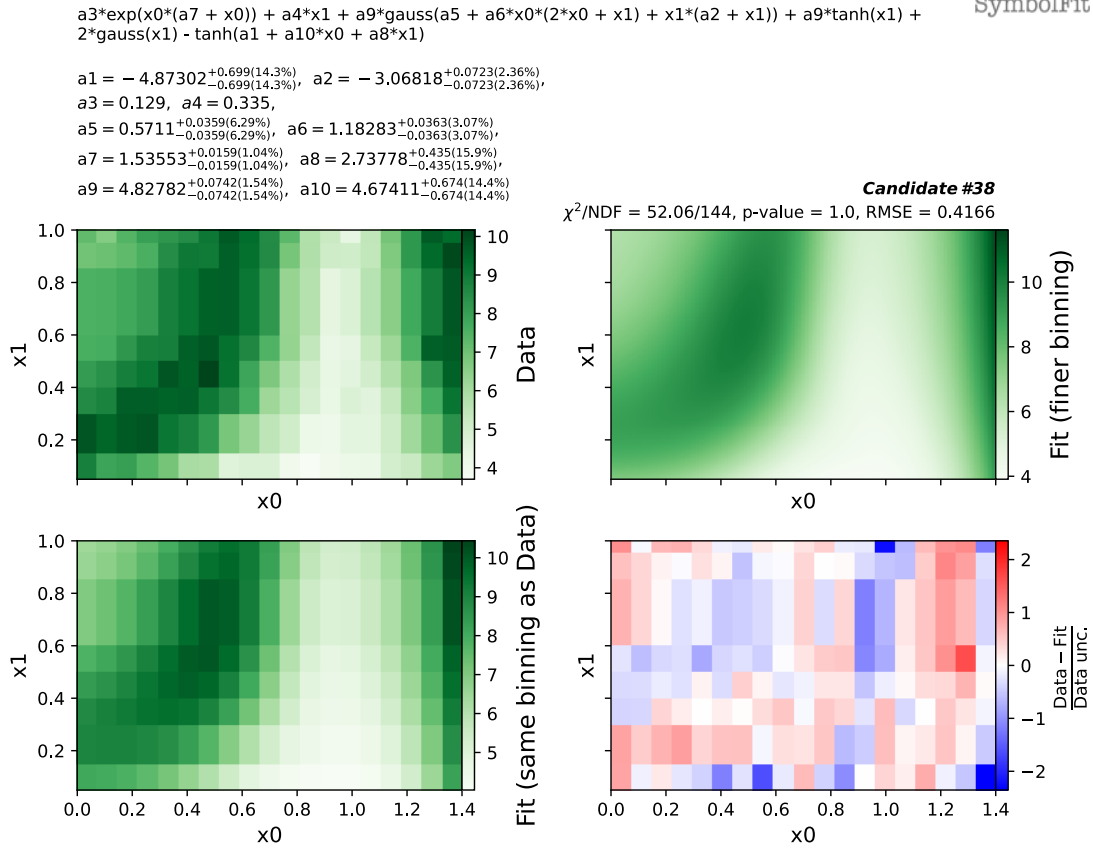


Figure 15: Candidate function #38 for Toy Dataset 3b (see Tab. 4). The parameterized form of this function is shown at the top of the figure, along with the best-fit values and associated uncertainties. Upper left: the binned data being fitted. Lower left: the candidate function plotted with the same binning as the fitting data. Upper right: the candidate function plotted with a finer binning. Lower right: the residual error, $\frac{\text{Data} - \text{Fit}}{\text{Data unc.}}$, in units of the data uncertainty.

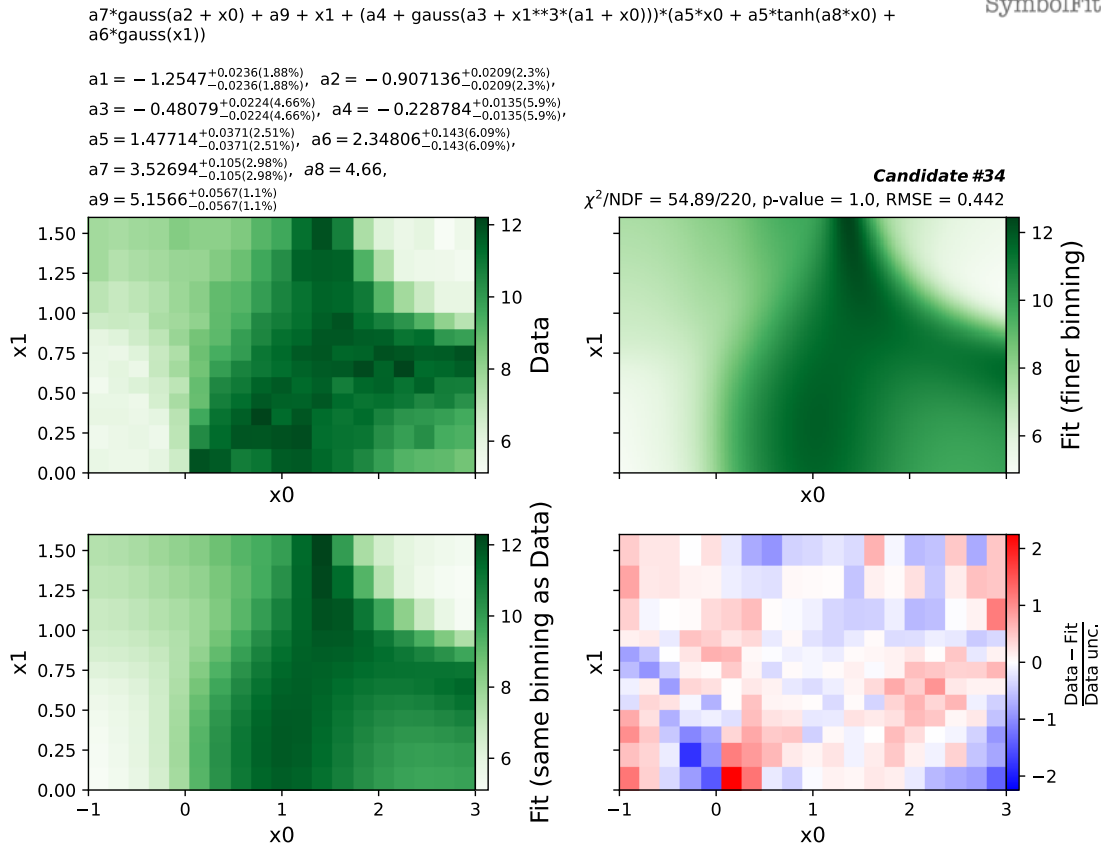


Figure 16: Candidate function #34 for Toy Dataset 3c (see Tab. 4). The parameterized form of this function is shown at the top of the figure, along with the best-fit values and associated uncertainties. Upper left: the binned data being fitted. Lower left: the candidate function plotted with the same binning as the fitting data. Upper right: the candidate function plotted with a finer binning. Lower right: the residual error, $\frac{\text{Data} - \text{Fit}}{\text{Uncertainty}}$, in units of the data uncertainty.

6. Summary

We have developed a framework called `SymbolFit` that automates parametric modeling without the need for a priori specification of a functional form to fit data. The framework employs symbolic regression to machine-search for functional forms and incorporates a re-optimization step to improve the candidate functions and provide uncertainty estimates. Due to the nature of genetic programming, each symbolic regression fit generates a batch of candidate functions with a variety of forms that can potentially model the data well. This offers flexibility, allowing users to select the most suitable candidates for whatever downstream tasks

Our primary focus is on applications in high-energy physics data analysis, specifically in signal and background modeling, as well as the derivation of smooth scale factors from binned data. We have demonstrated the effectiveness and efficiency of our framework using five real proton-proton collision datasets from new physics searches at the CERN LHC, as well as several toy datasets, including two-dimensional binned data.

Since the fit outputs in this approach are parametric closed-form functions, the resulting model representation is identical to that from traditional parametric modeling methods. This allows seamless integration with established downstream statistical tools such as `Combine` and `pyhf` for hypothesis testing without requiring separate treatment. Furthermore, the ease of generating a wide range of well-fitted functions within this framework facilitates flexible modeling, as the choice of functions can be treated as a source of systematic uncertainty using well-established techniques, such as the discrete profiling method.

We have also developed an API for the framework, designed for easy use by the high-energy physics community. This API automates the process of finding suitable functions for modeling binned data, significantly reducing manual effort. Our goal is to transform the approach to parametric modeling in high-energy physics experiments, moving away from traditional fitting methods that rely on manually determining empirical functions on a case-by-case basis, which is both time-consuming and prone to bias.

Acknowledgements

DR is supported by the U.S. Department of Energy (DOE), Office of Science, Office of High Energy Physics Early Career Research program under Award No. DE-SC0025324. SD is supported by the U.S. DOE, Office of Science, Office of High Energy Physics, under Award No. DE-SC0017647. JD is supported by the Research Corporation for Science Advancement (RCSA) under grant #CS-CSA-2023-109, Alfred P. Sloan Foundation under grant #FG-2023-20452, U.S. DOE, Office of Science, Office of High Energy Physics Early Career Research program under Award No. DE-SC0021187. JD, PH, and VL are supported by the U.S. National Science Foundation (NSF) Harnessing the Data Revolution (HDR) Institute for Accelerating AI Algorithms for Data Driven Discovery (A3D3) under Cooperative Agreement PHY-2117997. PH is also supported by the Institute for Artificial Intelligence and Fundamental Interactions (IAIFI), under the NSF grant #PHY-2019786. EL is supported by the U.S. DOE, Office of Science, Office of High Energy Physics, under Award No. DE-SC0007901.

References

- [1] La Cava, W. *et al.* Contemporary symbolic regression methods and their relative performance. In Vanschoren, J. & Yeung, S. (eds.) *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks*, vol. 1 (2021). URL https://datasets-benchmarks-proceedings.neurips.cc/paper_files/paper/2021/file/c0c7c76d30bd3dcaefc96f40275bdc0a-Paper-round1.pdf.
- [2] Koza, J. Genetic programming as a means for programming computers by natural selection. *Statistics and Computing* **4**, 87–112 (1994).
- [3] Schmidt, M. & Lipson, H. Distilling free-form natural laws from experimental data. *Science* **324**, 81–85 (2009). URL <https://doi.org/10.1126/science.1165893>.

- [4] Stephens, T. Genetic programming in Python, with a scikit-learn inspired API: gplearn (2016). URL <https://gplearn.readthedocs.io/en/stable/>.
- [5] Burlacu, B., Kronberger, G. & Kommenda, M. Operon C++: An efficient genetic programming framework for symbolic regression. In *Proceedings of the 2020 Genetic and Evolutionary Computation Conference Companion*, GECCO '20, 1562–1570 (Association for Computing Machinery, New York, NY, USA, 2020). URL <https://doi.org/10.1145/3377929.3398099>.
- [6] Virgolin, M., Alderliesten, T., Witteveen, C. & Bosman, P. A. N. Improving model-based genetic programming for symbolic regression of small expressions. *Evolutionary Computation* **29**, 211–237 (2021). URL https://doi.org/10.1162/2Fevco_a_00278.
- [7] Cranmer, M. Interpretable Machine Learning for Science with PySR and SymbolicRegression.jl (2023). 2305.01582.
- [8] Hayrapetyan, A. *et al.* Search for Narrow Trijet Resonances in Proton-Proton Collisions at $\sqrt{s}=13$ TeV. *Phys. Rev. Lett.* **133**, 011801 (2024). 2310.14023.
- [9] Fisher, R. A. On the interpretation of χ^2 from contingency tables, and the calculation of p . *Journal of the Royal Statistical Society* **85**, 87–94 (1922). URL <http://www.jstor.org/stable/2340521>.
- [10] Aad, G. *et al.* Observation of a new particle in the search for the Standard Model Higgs boson with the ATLAS detector at the LHC. *Phys. Lett. B* **716**, 1–29 (2012). 1207.7214.
- [11] Chatrchyan, S. *et al.* Observation of a New Boson at a Mass of 125 GeV with the CMS Experiment at the LHC. *Phys. Lett. B* **716**, 30–61 (2012). 1207.7235.
- [12] Chatrchyan, S. *et al.* Observation of a New Boson with Mass Near 125 GeV in pp Collisions at $\sqrt{s} = 7$ and 8 TeV. *JHEP* **06**, 081 (2013). 1303.4571.
- [13] Sirunyan, A. M. *et al.* Search for high mass dijet resonances with a new background prediction method in proton-proton collisions at $\sqrt{s} = 13$ TeV. *JHEP* **05**, 033 (2020). 1911.03947.
- [14] Tumasyan, A. *et al.* Search for resonant and nonresonant production of pairs of dijet resonances in proton-proton collisions at $\sqrt{s} = 13$ TeV. *JHEP* **07**, 161 (2023). 2206.09997.
- [15] Hayrapetyan, A. *et al.* Search for new physics in high-mass diphoton events from proton-proton collisions at $\sqrt{s} = 13$ TeV. *JHEP* **08**, 215 (2024). 2405.09320.
- [16] Hayrapetyan, A. *et al.* Search for a high-mass dimuon resonance produced in association with b quark jets at $\sqrt{s} = 13$ TeV. *JHEP* **10**, 043 (2023). 2307.08708.
- [17] Khachatryan, V. *et al.* Jet energy scale and resolution in the CMS experiment in pp collisions at 8 TeV. *JINST* **12**, P02014 (2017). 1607.03663.
- [18] Sirunyan, A. M. *et al.* Identification of heavy-flavour jets with the CMS detector in pp collisions at 13 TeV. *JINST* **13**, P05011 (2018). 1712.07158.
- [19] Sirunyan, A. M. *et al.* Performance of reconstruction and identification of τ leptons decaying to hadrons and ν_τ in pp collisions at $\sqrt{s} = 13$ TeV. *JINST* **13**, P10005 (2018). 1809.02816.
- [20] Sirunyan, A. M. *et al.* Inclusive search for highly boosted Higgs bosons decaying to bottom quark-antiquark pairs in proton-proton collisions at $\sqrt{s} = 13$ TeV. *JHEP* **12**, 085 (2020). 2006.13251.
- [21] Frate, M., Cranmer, K., Kalia, S., Vandenberg-Rodes, A. & Whiteson, D. Modeling Smooth Backgrounds and Generic Localized Signals with Gaussian Processes (2017). 1709.05681.
- [22] Gandrakota, A., Lath, A., Morozov, A. V. & Murthy, S. Model selection and signal extraction using Gaussian Process regression. *JHEP* **02**, 230 (2023). 2202.05856.
- [23] Xu, R. A Measurement of Boosted Dibosons with Gaussian Process Background Modeling at the ATLAS Detector (2024). URL <https://cds.cern.ch/record/2901208>. Presented 21 May 2024.
- [24] Swiler, L. P., Gulian, M., Frankel, A. L., Safta, C. & Jakeman, J. D. A survey of constrained gaussian process regression: Approaches and implementation challenges. *Journal of Machine Learning for Modeling and Computing* **1**, 119–156 (2020). URL <http://dx.doi.org/10.1615/JMachLearnModelComput.2020035155>.
- [25] Neville, M., Stensitzki, T., Allen, D. B. & Ingargiola, A. LMFIT: Non-Linear Least-Square Minimization and Curve-Fitting for Python (2015). URL <https://doi.org/10.5281/zenodo.11813>.
- [26] Hayrapetyan, A. *et al.* The CMS Statistical Analysis and Combination Tool: COMBINE (2024). 2404.06614.
- [27] Heinrich, L., Feickert, M. & Stark, G. pyhf: v0.7.6. URL <https://doi.org/10.5281/zenodo.1169739>. <https://github.com/scikit-hep/pyhf/releases/tag/v0.7.6>.
- [28] Heinrich, L., Feickert, M., Stark, G. & Cranmer, K. pyhf: pure-python implementation of histfactory statistical models. *Journal of Open Source Software* **6**, 2823 (2021). URL <https://doi.org/10.21105/joss.02823>.
- [29] Dauncey, P. D., Kenzie, M., Wardle, N. & Davies, G. J. Handling uncertainties in background shapes: the discrete profiling method. *JINST* **10**, P04015 (2015). 1408.6865.
- [30] CMS Collaboration. Search for high mass dijet resonances with a new background prediction method in proton-proton collisions at $\sqrt{s} = 13$ TeV. HEPData (collection) (2019). <https://doi.org/10.17182/hepdata.91059>.
- [31] CMS Collaboration. Search for new physics in high-mass diphoton events from proton-proton collisions at $\sqrt{s} = 13$ TeV. HEPData (collection) (2024). <https://doi.org/10.17182/hepdata.150677>.
- [32] CMS Collaboration. Search for narrow trijet resonances in proton-proton collisions at $\sqrt{s} = 13$ TeV. HEPData (collection) (2023). <https://doi.org/10.17182/hepdata.144165>.
- [33] CMS Collaboration. Search for resonant and nonresonant production of pairs of dijet resonances in proton-proton collisions at $\sqrt{s} = 13$ TeV. HEPData (collection) (2022). <https://doi.org/10.17182/hepdata.130817>.

- [34] CMS Collaboration. Search for a high-mass dimuon resonance produced in association with b quark jets at $\sqrt{s}=13$ TeV. HEPData (collection) (2023). <https://doi.org/10.17182/hepdata.141455>.

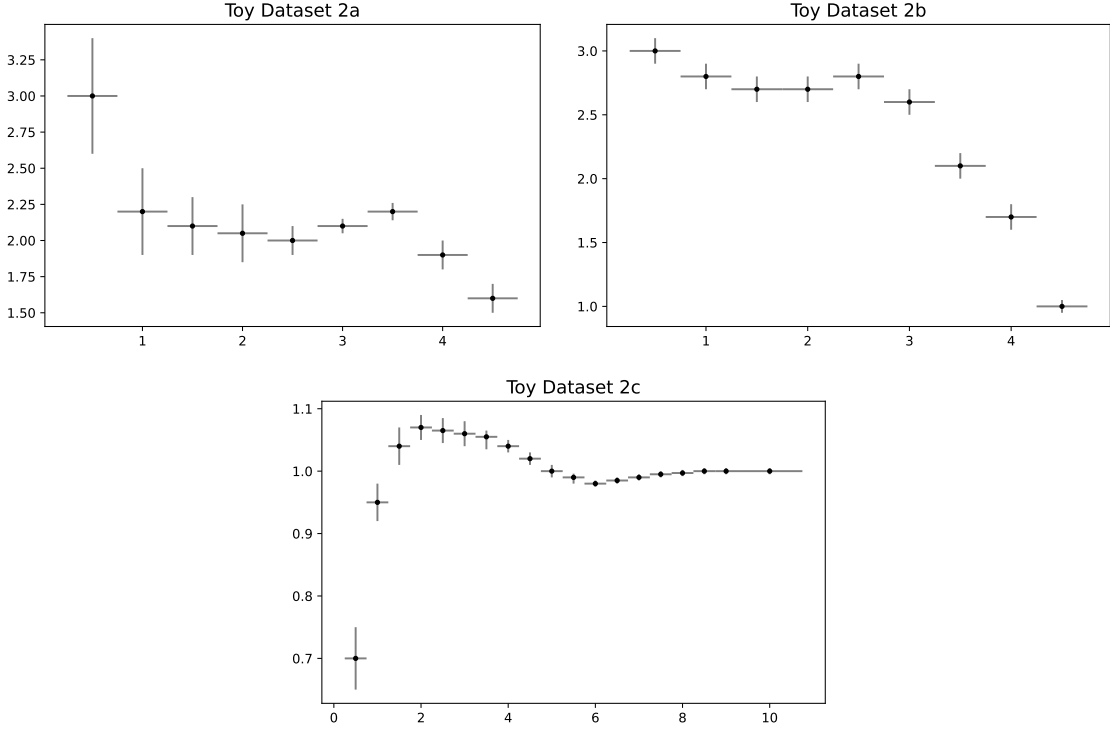


Figure A1: Toy Dataset 2: three 1D binned sub-datasets manually generated without reference to an underlying function.

Appendix A. More examples

Appendix A.1. Toy dataset 2 (1D) [arbitrary shapes]

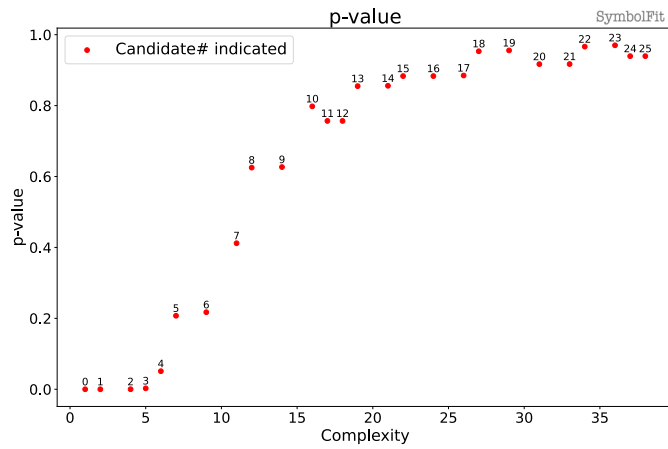
In Toy Dataset 2, we consider three 1D binned sub-datasets, labeled 2a, 2b, and 2c, as shown in Fig. A1. These datasets are manually generated without reference to an underlying function and are used to demonstrate applications such as deriving smooth scale factors from binned data.

We use the same PySR configuration applied to the five LHC datasets, as shown in List. 2, to fit these 1D binned datasets. For each sub-dataset, a single run of SymbolFit is performed to generate a batch of candidate functions. Fig. A2 shows the p-value plotted against function complexity. Several candidate functions for each sub-datasets are selected and listed in Tab. A1.

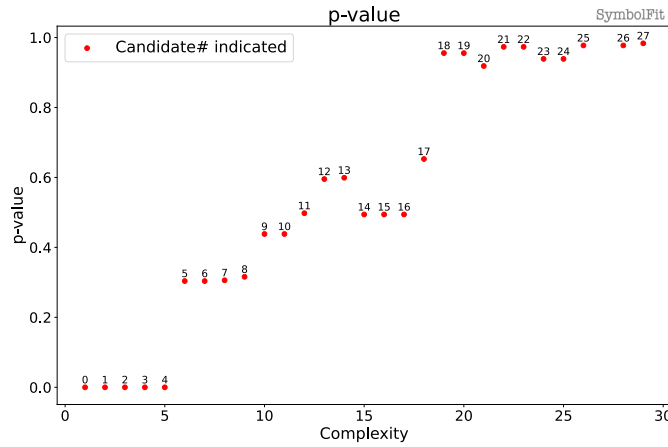
Fig. A3 shows candidate functions #13, #21, and #21 with uncertainty coverage, for sub-datasets 2a, 2b, and 2c, respectively.

Table A1: Examples candidate functions for Toy Dataset 2 are listed. The example candidate functions—#13 for Toy Dataset 2a, #21 for 2b, and #21 for 2c—are plotted in Fig. A3. Numerical values are rounded to three significant figures for display purposes.

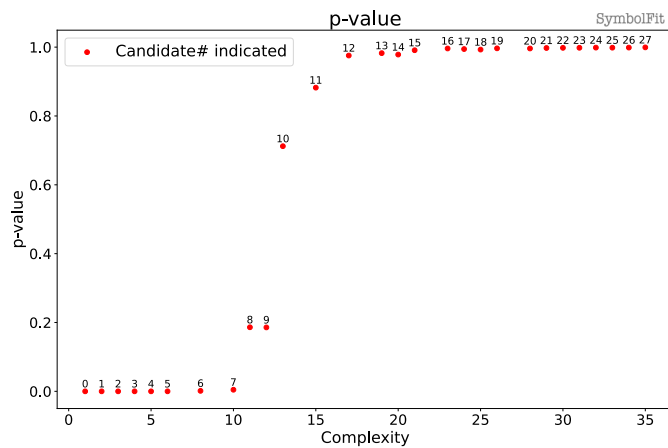
Toy Dataset 2a					
Complexity	Candidate function (after ROF)	# param.	χ^2/NDF (before ROF)	χ^2/NDF (after ROF)	p-value (after ROF)
12 (#8)	$0.667 \times 0.991^{\exp(x)}x + 1.44/x$	3	5.273 / 6 = 0.8789	4.382 / 6 = 0.7303	0.6251
19 (#13)	$0.697 \times 0.991^{\exp(x)}x + 0.697^{\exp(x)}x + 1.16/x$	3	2.863 / 6 = 0.4772	2.619 / 6 = 0.4365	0.8549
38 (#25)	$0.907^{0.0993 \exp(x)}x(0.0609 + \tanh(0.596^{0.632/x})) + 0.948/x + 0.948^{\exp(1.63x)} - 0.115$	3	1.776 / 6 = 0.2961	1.771 / 6 = 0.2951	0.9395
Toy Dataset 2b					
Complexity	Candidate function (after ROF)	# param.	χ^2/NDF (before ROF)	χ^2/NDF (after ROF)	p-value (after ROF)
14 (#13)	$-0.0219 \exp(x) + 1.96 + \tanh(x)^{-0.586+x}$	2	5.613 / 7 = 0.8018	5.501 / 7 = 0.7858	0.5991
19 (#18)	$-0.0161 \exp(x) + 0.177 + 3.02 \times 0.825^x + \tanh(0.177x^x)$	4	2.092 / 5 = 0.4184	1.084 / 5 = 0.2169	0.9555
22 (#21)	$-0.0149 \exp(x) - 0.00659 + 4.11 \tanh(0.787^x) + \tanh(0.16x^x)$	4	0.857 / 5 = 0.1714	0.8561 / 5 = 0.1712	0.9733
Toy Dataset 2c					
Complexity	Candidate function (after ROF)	# param.	χ^2/NDF (before ROF)	χ^2/NDF (after ROF)	p-value (after ROF)
13 (#10)	$0.626x \exp(-1.38^x) + \tanh(x^{0.502})$	3	16.75 / 16 = 1.047	12.45 / 16 = 0.7784	0.7122
15 (#11)	$x/(2x + \exp(1.41^x)) + \tanh(x^{0.505})$	2	10.67 / 17 = 0.6276	10.48 / 17 = 0.6165	0.8823
29 (#21)	$(1.02x + 1.02 \tanh(x^2)) \tanh(x)/(2.58x^{1.39} + x + \exp(1.39^x)) + 1.02 \tanh(x^{0.389})$	4	7.767 / 15 = 0.5178	4.079 / 15 = 0.2719	0.9975



(a) Toy Dataset 2a.

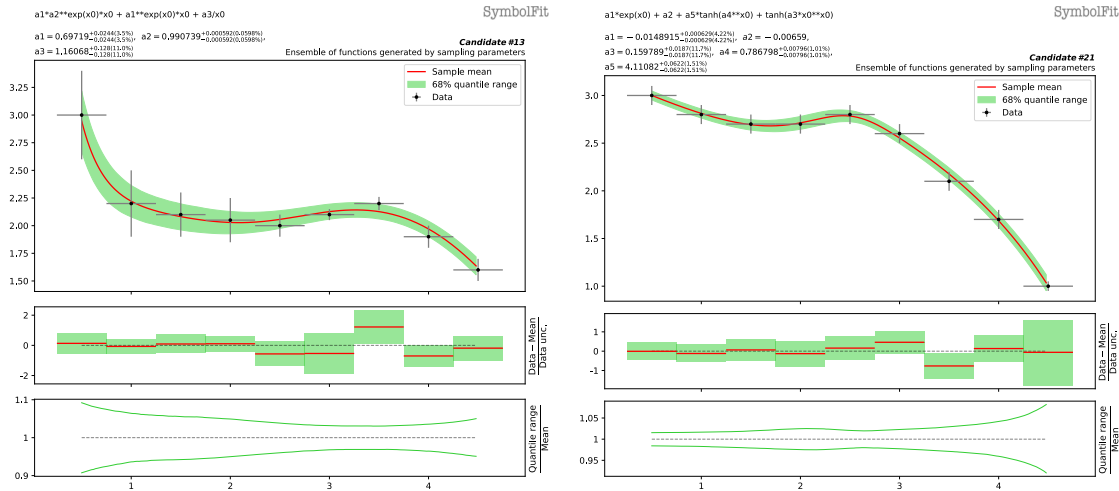


(b) Toy Dataset 2b.

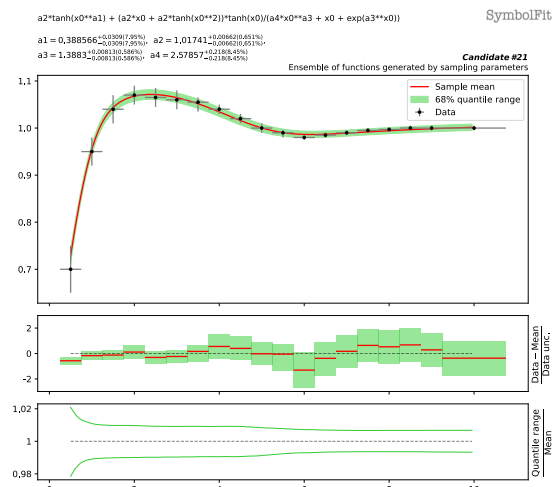


(c) Toy Dataset 2c.

Figure A2: p-value vs. function complexity. A total of 26, 28, and 28 candidate functions (labeled #0–#25, #0–#27, and #0–#27) were obtained from a single fit on Toy Dataset 2a, 2b, and 2c, respectively.



(a) Candidate function #13 for Toy Dataset 2a. (b) Candidate function #21 for Toy Dataset 2b.



(c) Candidate function #21 for Toy Dataset 2c.

Figure A3: Example candidate functions for Toy Dataset 2 (see Tab. A1). To visualize the total uncertainty coverage of each candidate function, the green band in each subfigure represents the 68% quantile range of functions obtained by sampling parameters, taking into account the best-fit values and the covariance matrix within a multidimensional normal distribution. The red line denotes the mean of the function ensemble. At the top of each subfigure, the candidate function and the fitted parameters are shown. The middle panel shows the weighted residual error: $\frac{\text{Data} - \text{Mean}}{\text{Data unc.}}$. The bottom panel shows the ratio of the 68% quantile range to the mean.

Table A2: The candidate functions are obtained from three fits using different random seeds, fitted to the pseudodata of the diphoton spectrum with the (injected) signal region blinded. The fits were performed on a scaled dataset (to enhance fit stability and prevent numerical overflow), and the functions can be transformed back to describe the original spectrum using the transformation: $x \rightarrow 0.000287(x - 503)$. These functions are plotted and compared with the blinded pseudodata in Fig. A5. Numerical values are rounded to three significant figures for display purposes.

	Candidate function (after ROF)	# param.	χ^2/NDF (before ROF)	χ^2/NDF (after ROF)	p-value (after ROF)
SR model 1	$33.3(0.017 + 0.177x)^{6.93x}$	3	47.12 / 136 = 0.3465	46.83 / 136 = 0.3443	1.0
SR model 2	$35.0 \exp(-16.0x^{0.834})$	3	63.44 / 136 = 0.4665	53.37 / 136 = 0.3925	1.0
SR model 3	$0.0226 \exp(0.975x)(0.158 + x)^{-3.96}$	3	48.94 / 136 = 0.3598	47.22 / 136 = 0.3472	1.0

Appendix A.2. CMS high-mass diphoton dataset (1D) [background modeling]

CMS performed a search for high-mass diphoton resonances using proton-proton collision data at a center-of-mass energy of $\sqrt{s} = 13$ TeV and reported no significant deviations from the Standard Model prediction [15]. The dataset for the diphoton spectrum is publicly available on HEPDATA at Ref. [31]. In the analysis, CMS considered four empirical functions to model the background contribution in the distribution of the diphoton invariant mass, $m_{\gamma\gamma}$, and one of them is:

$$f(x) = p_0 x^{p_1 + p_2 \log(x)}, \quad (\text{A.1})$$

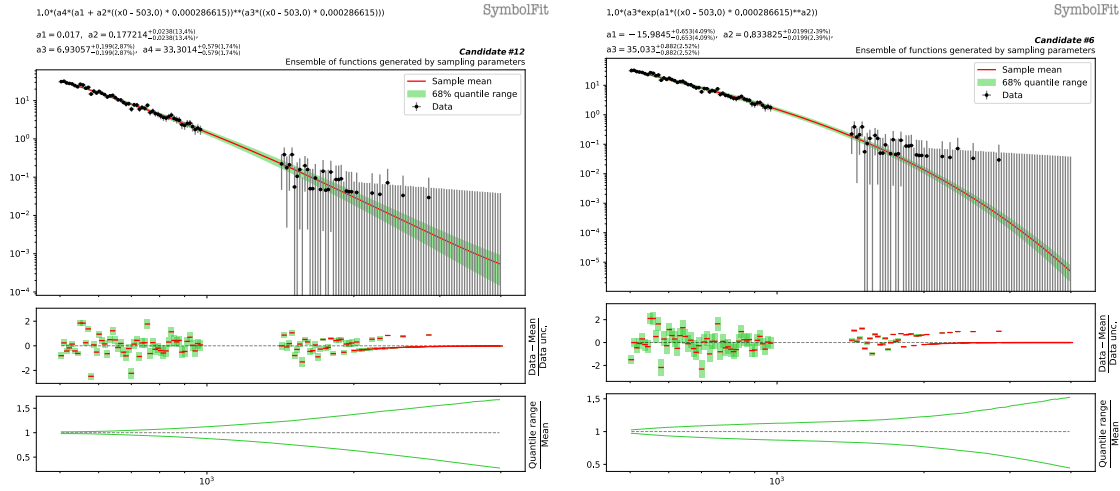
where $x = m_{\gamma\gamma}/\sqrt{s}$ is dimensionless and $p_{\{0,1,2\}}$ are free parameters.

We perform the same experiments conducted on the dijet dataset, as detailed in Sec. 5.2. Starting from the original diphoton spectrum, we generate pseudodata by injecting a perturbed Gaussian signal centered at $m_{\gamma\gamma} = 1320$ GeV (s_1), with a width of 74 GeV ($2s_2$) and a signal strength of $s_0 = 15$. To model the background, we blind the signal region by masking the $m_{\gamma\gamma}$ bins between 980 and 1400 GeV in the pseudodata and perform the fits.

Three `SymbolFit` runs using different random seeds are carried out, applying the same `PySR` configuration as used for the dijet dataset (see List. 2), except that the maximum complexity is set at 20 instead of 80, since the diphoton distribution shape is less complex. Tab. A2 lists the three SR models, each obtained from a fit initialized with a different random seed. The χ^2/NDF scores improve significantly after the ROF step compared to the original functions returned by `PySR`. The three background models fit the blinded pseudodata well, as shown in Fig. A4 for the total uncertainty coverage and Fig. A5 for a comparison with the empirical model used by CMS.

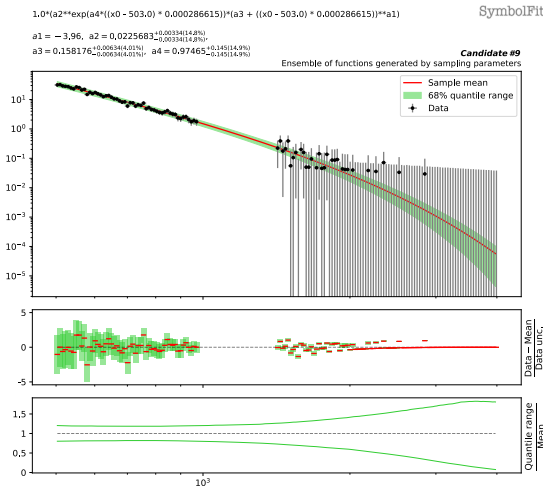
Next, we unblind the pseudodata and perform b-only fits and s+b fits on the full pseudodata spectrum. These results are shown in Fig. A6. In all three SR models, as well as the CMS empirical model, the excess of events over the background around the injected signal location observed in the b-only fits is reduced in the s+b fits, demonstrating that the models are sensitive to the injected signal. Tab. A3 lists the χ^2/NDF scores for each model, showing the fit performance in response to the presence of the injected signal.

To assess whether the SR models can accurately extract the injected signals, we generate multiple sets of pseudodata by injecting Gaussian signals with different mean values ranging from 1080 to 1120 GeV and varying signal strengths between 5 and 50. We then perform the s+b fits to extract the corresponding signal parameters. Fig. A7 shows the extracted signal parameters plotted against their injected values. All three SR models are capable of extracting the correct signal parameter values within reasonable uncertainties and are comparable to the empirical model used by CMS.



(a) SR model 1.

(b) SR model 2.



(c) SR model 3.

Figure A4: The three SR models fitted to the pseudodata of the diphoton spectrum with the signal region blinded (see Tab. A2). To visualize the total uncertainty coverage of each candidate function, the green band in each subfigure represents the 68% quantile range of functions obtained by sampling parameters, taking into account the best-fit values and the covariance matrix within a multidimensional normal distribution. The red line denotes the mean of the function ensemble. At the top of each subfigure, the candidate function and the fitted parameters are shown. The middle panel shows the weighted residual error: $\frac{\text{Data} - \text{Mean}}{\text{Data unc.}}$. The bottom panel shows the ratio of the 68% quantile range to the mean.

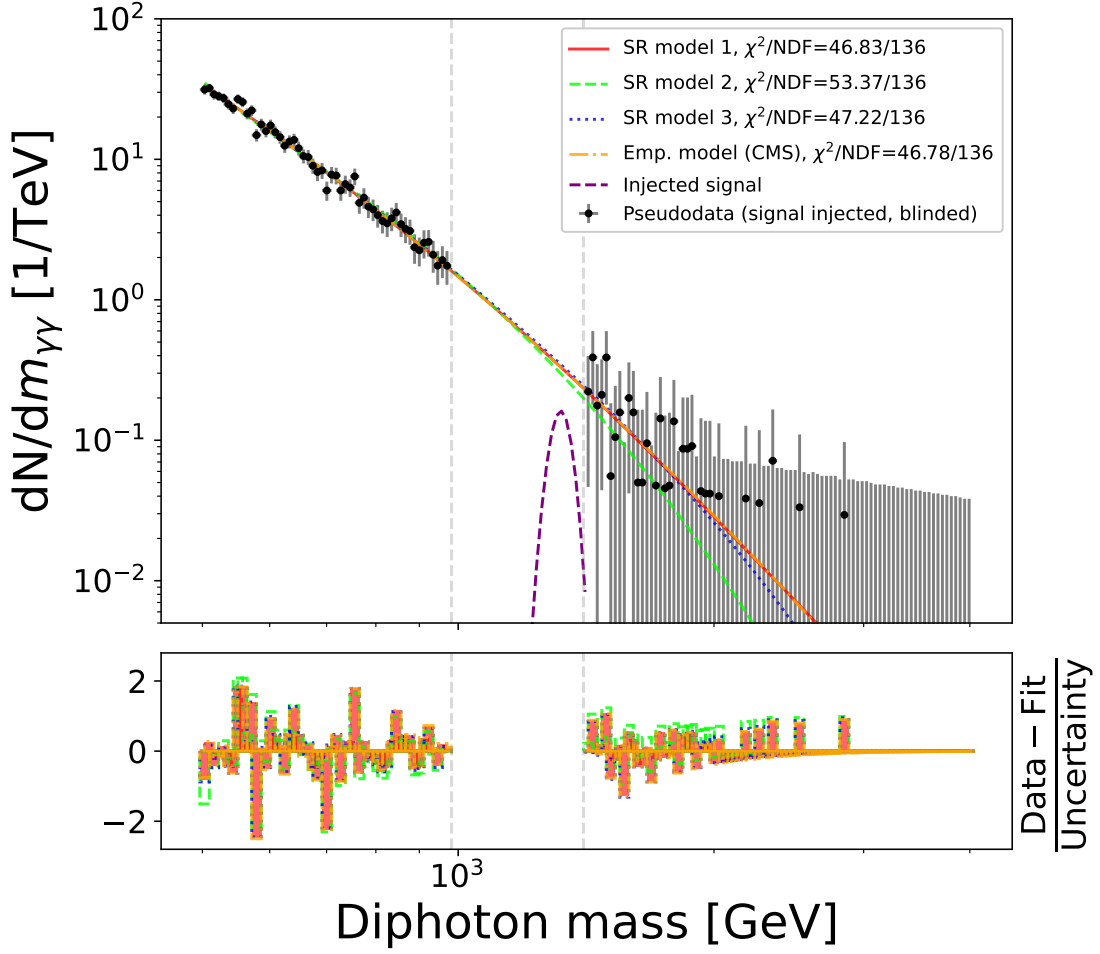


Figure A5: Pseudodata of the diphoton spectrum with the injected signal shown in the blinded signal region. The three SR models (see Tab. A2) are compared against the empirical model used by CMS. The lower panel shows the residual error per bin, measured in units of the data uncertainty.

Table A3: Comparison of the χ^2/NDF scores from three types of fits to the diphoton dataset: the b-only fits to the blinded pseudodata, b-only fits to the unblinded pseudodata, and s+b fits to the unblinded pseudodata. The background models used for the fits are listed in Tab. A2, and the fits are shown in Fig. A5 (blinded) and Fig. A6 (unblinded).

	χ^2/NDF (b-only, blinded)	χ^2/NDF (b-only, unblinded)	χ^2/NDF (s+b, unblinded)
SR model 1	46.83 / 136 = 0.3443	70.27 / 166 = 0.4233	62.34 / 164 = 0.3801
SR model 2	53.37 / 136 = 0.3924	79.94 / 166 = 0.4816	69.62 / 164 = 0.4245
SR model 3	47.22 / 136 = 0.3472	70.59 / 166 = 0.4252	63.26 / 164 = 0.3857
Emp. model (CMS)	46.78 / 136 = 0.344	70.02 / 166 = 0.4218	62.2 / 164 = 0.3793

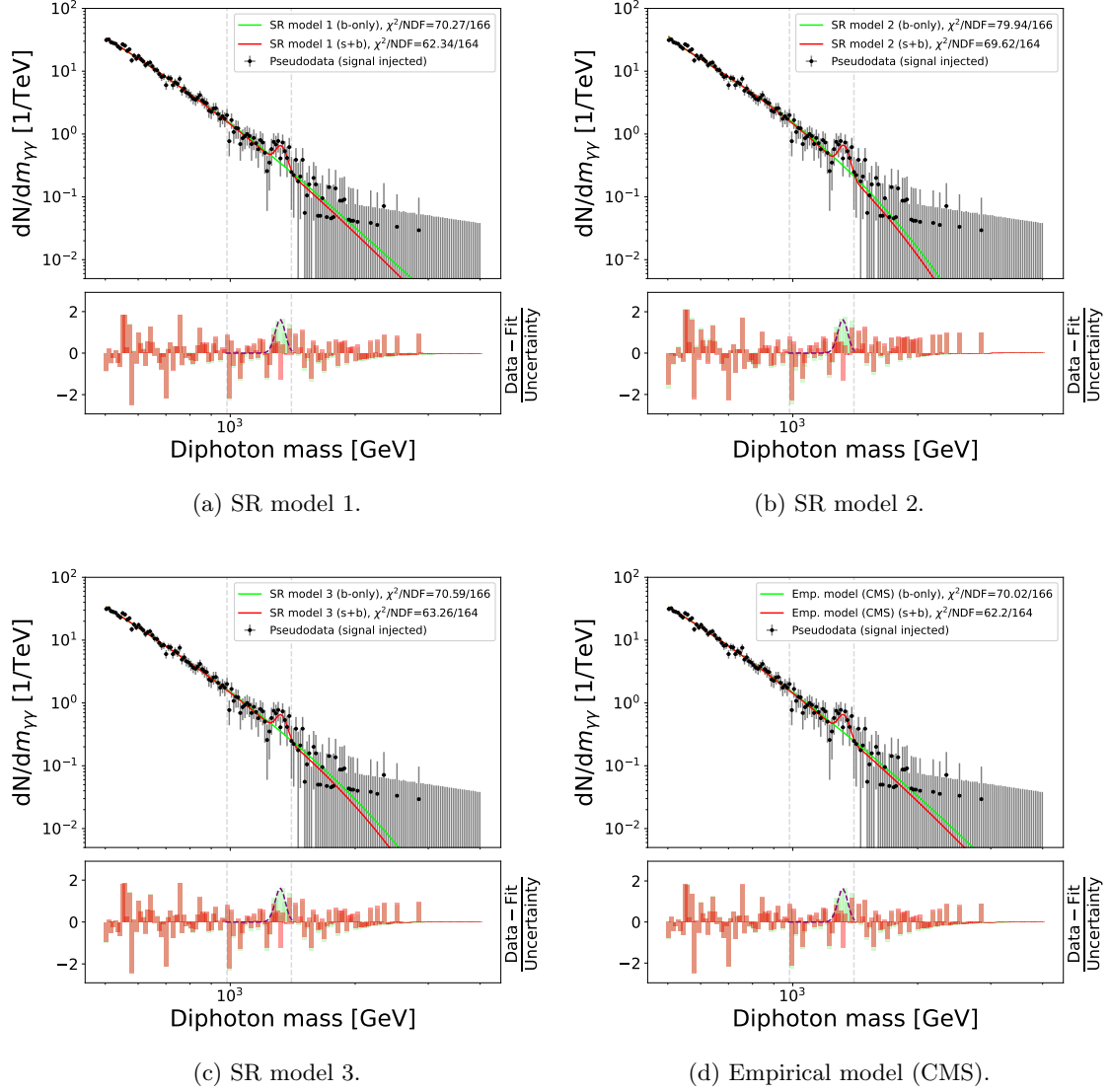
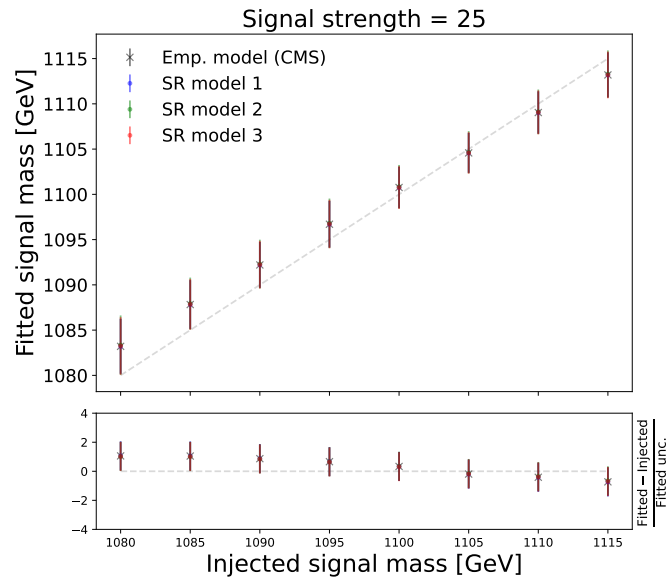
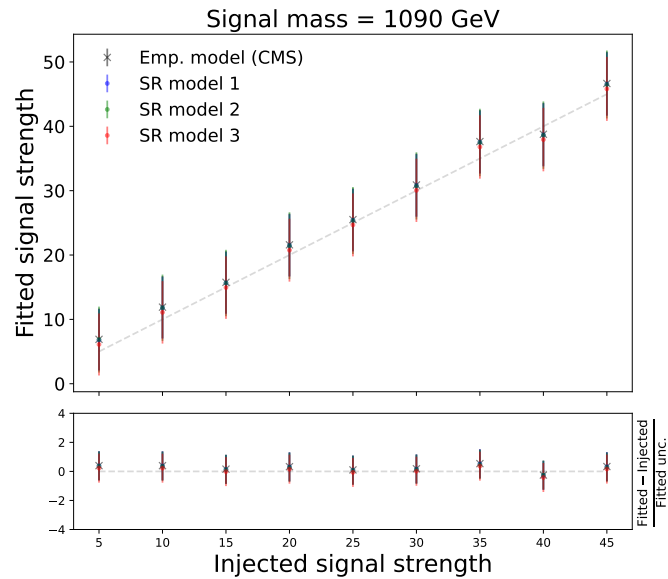


Figure A6: Comparison of the b-only fits and the s+b fits to the unblinded pseudodata of the diphoton spectrum. The lower panel shows the residual error per bin, measured in units of the data uncertainty. The shape of the injected signal is also shown.



(a) Fitted vs. injected signal mass at a specified signal strength value.



(b) Fitted vs. injected signal strength at a specified signal mass value.

Figure A7: Fitted values vs. the true values of parameters of the injected signal in the diphoton dataset. The bottom panels show the residual error in units of the fitted uncertainty.

Table A4: The candidate functions are obtained from three fits using different random seeds, fitted to the pseudodata of the trijet spectrum with the (injected) signal region blinded. The fits were performed on a scaled dataset (to enhance fit stability and prevent numerical overflow), and the functions can be transformed back to describe the original spectrum using the transformation: $f(x) \rightarrow 38458 \times f(0.000184(x - 1790))$. These functions are plotted and compared with the blinded pseudodata in Fig. A9. Numerical values are rounded to three significant figures for display purposes.

	Candidate function (after ROF)	# param.	χ^2/NDF (before ROF)	χ^2/NDF (after ROF)	p-value (after ROF)
SR model 1	$(1.08 \times 10^{-5})^{\tanh(x)} / ((0.165 + x) \times \exp(x^2(-1.96 + 4x))^{\tanh(1.17x^2)})$	3	50.46 / 26 = 1.941	49.04 / 26 = 1.886	0.00408
SR model 2	$\exp(x(-10.8 + x)) / (-0.261x \tanh(x \times (-10.9 + x)) + 0.165 + \tanh(x))$	4	39.49 / 25 = 1.58	33.15 / 25 = 1.326	0.1273
SR model 3	$0.0554^{-0.622+4.03x} (0.568 \tanh(2x) + (0.00302 \exp(x)/(1.92 + x))^x)$	4	38.4 / 25 = 1.536	37.31 / 25 = 1.492	0.05395

Appendix A.3. CMS trijet dataset (1D) [background modeling]

CMS performed a search for high-mass trijet resonances using proton-proton collision data at a center-of-mass energy of $\sqrt{s} = 13$ TeV and reported no significant deviations from the Standard Model prediction [8]. The dataset for the trijet spectrum is publicly available on HEPDATA at Ref. [32]. In the analysis, CMS considered four empirical functions to model the background contribution in the distribution of the trijet invariant mass, m_{jjj} , and one of them is:

$$f(x) = \frac{p_0(1-x)^{p_1}}{x^{p_2+p_3} \log(x)}, \quad (\text{A.2})$$

where $x = m_{\text{jjj}}/\sqrt{s}$ is dimensionless and $p_{\{0,1,2,3\}}$ are free parameters. Eq. A.3 corresponds to Eq. 1 with $N = 3$ determined by an F-test.

We perform the same experiments conducted on the dijet dataset, as detailed in Sec. 5.2. Starting from the original trijet spectrum, we generate pseudodata by injecting a perturbed Gaussian signal centered at $m_{\text{jjj}} = 4000$ GeV (s_1) with a width of 400 GeV ($2s_2$) and a signal strength of $s_0 = 50000$. To model the background, we blind the signal region by masking the m_{jjj} bins between 3000 and 5000 GeV in the pseudodata and perform the fits.

Three `SymbolFit` runs using different random seeds are carried out, applying the same `PySR` configuration as used for the dijet dataset (see List. 2). Tab. A4 lists the three SR models, each obtained from a fit initialized with a different random seed. The χ^2/NDF scores improve significantly after the ROF step compared to the original functions returned by `PySR`. The three background models fit the blinded pseudodata well, as shown in Fig. A8 for the total uncertainty coverage and Fig. A9 for a comparison with the empirical model used by CMS.

Next, we unblind the pseudodata and perform b-only fits and s+b fits on the full pseudodata spectrum. These results are shown in Fig. A10. In all three SR models, as well as the CMS empirical model, the excess of events over the background around the injected signal location observed in the b-only fits is reduced in the s+b fits, demonstrating that the models are sensitive to the injected signal. Tab. A5 lists the χ^2/NDF scores for each model, showing the fit performance in response to the presence of the injected signal.

To assess whether the SR models can accurately extract the injected signals, we generate multiple sets of pseudodata by injecting Gaussian signals with different mean values ranging from 3600 to 4500 GeV and varying signal strength between 25000 and 100000. We then perform the s+b

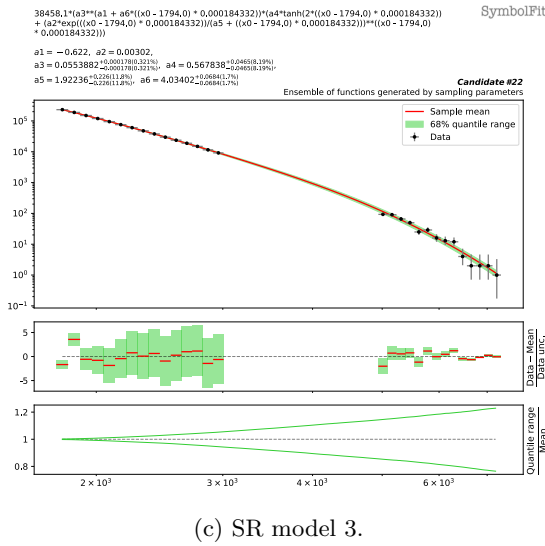
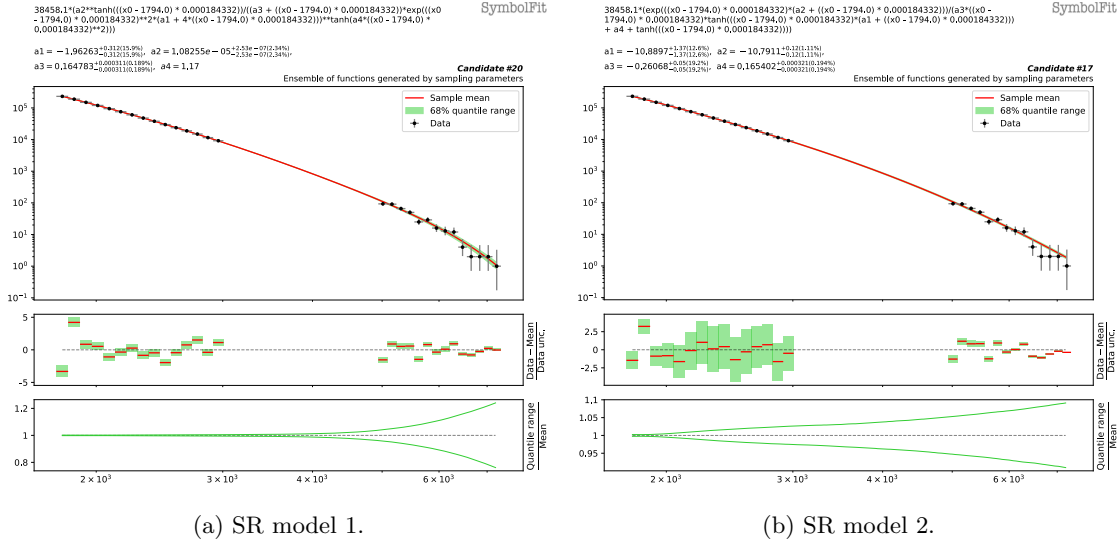


Figure A8: The three SR models fitted to the pseudodata of the trijet spectrum with the signal region blinded (see Tab. A4). To visualize the total uncertainty coverage of each candidate function, the green band in each subfigure represents the 68% quantile range of functions obtained by sampling parameters, taking into account the best-fit values and the covariance matrix within a multidimensional normal distribution. The red line denotes the mean of the function ensemble. At the top of each subfigure, the candidate function and the fitted parameters are shown. The middle panel shows the weighted residual error: $\frac{\text{Data} - \text{Mean}}{\text{Data unc.}}$. The bottom panel shows the ratio of the 68% quantile range to the mean.

fits to extract the corresponding signal parameters. Fig. A11 shows the extracted signal parameters plotted against their injected values. All three SR models are capable of extracting the correct signal parameter values within reasonable uncertainties and are comparable to the empirical model used by CMS.

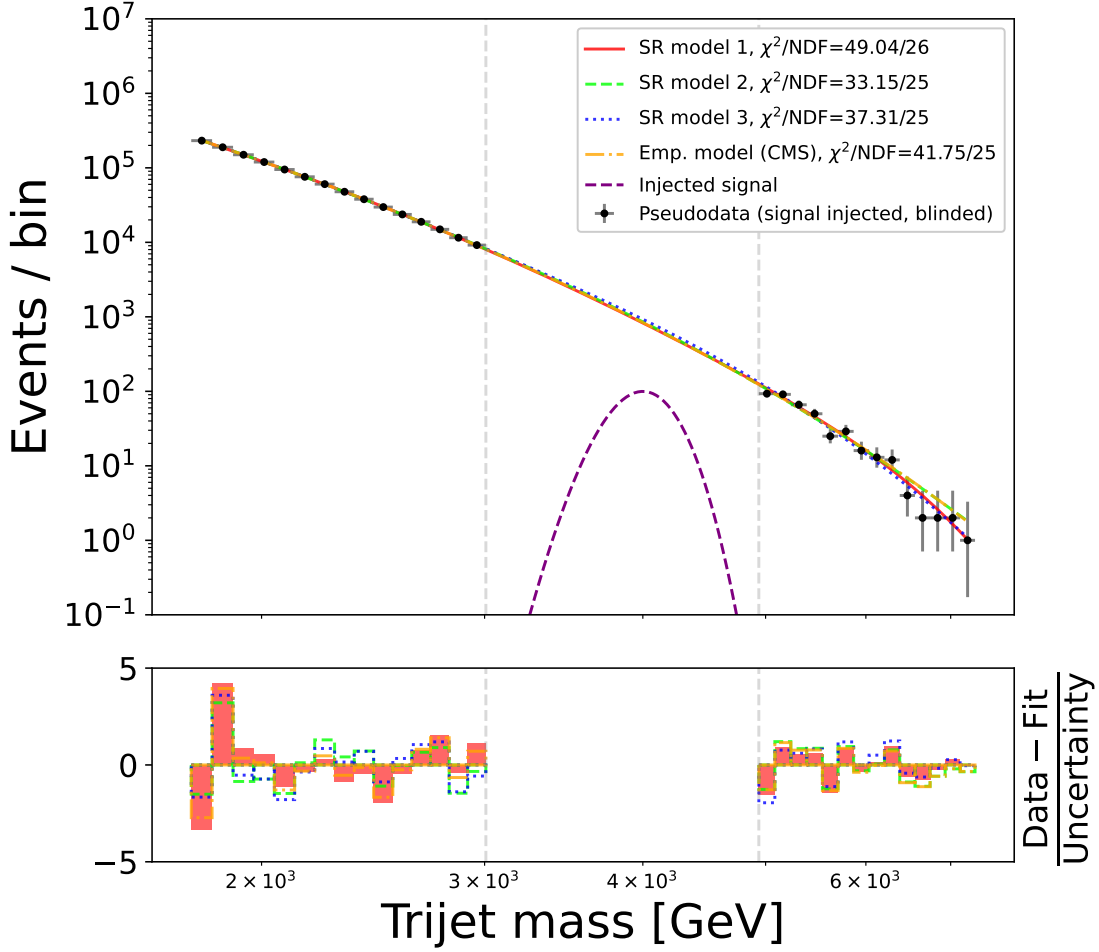


Figure A9: Pseudodata of the trijet spectrum with the injected signal shown in the blinded signal region. The three SR models (see Tab. A4) are compared against the empirical model used by CMS. The lower panel shows the residual error per bin, measured in units of the data uncertainty.

Table A5: Comparison of the χ^2/NDF scores from three types of fits to the trijet dataset: the b-only fits to the blinded pseudodata, b-only fits to the unblinded pseudodata, and s+b fits to the unblinded pseudodata. The background models used for the fits are listed in Tab. A4, and the fits are shown in Fig. A9 (blinded) and Fig. A10 (unblinded).

	χ^2/NDF (b-only, blinded)	χ^2/NDF (b-only, unblinded)	χ^2/NDF (s+b, unblinded)
SR model 1	49.04 / 26 = 1.886	181.3 / 42 = 4.317	93.04 / 40 = 2.326
SR model 2	33.15 / 25 = 1.326	92.92 / 41 = 2.266	48.07 / 39 = 1.233
SR model 3	37.31 / 25 = 1.492	74.5 / 41 = 1.817	54.23 / 39 = 1.391
Emp. model (CMS)	41.75 / 25 = 1.67	117.6 / 41 = 2.868	65.69 / 39 = 1.684

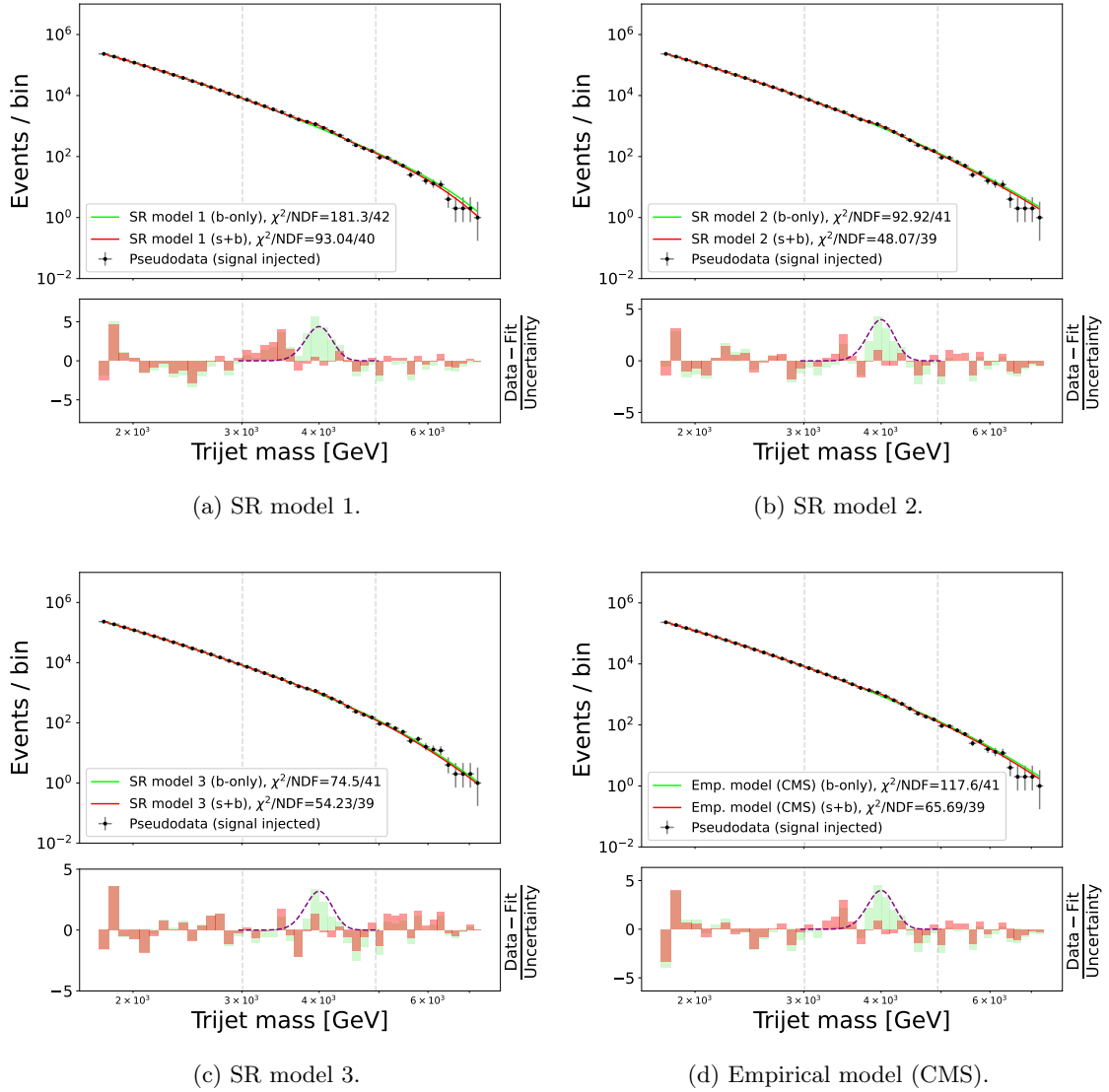
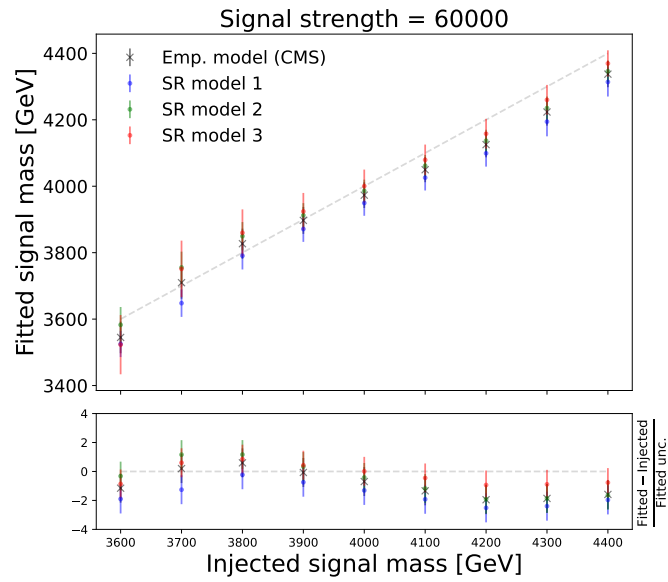
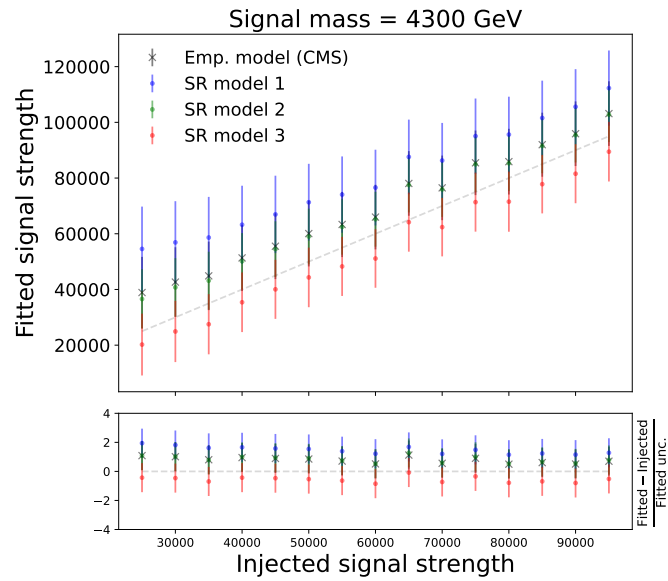


Figure A10: Comparison of the b-only fits and the s+b fits to the unblinded pseudodata. The lower panel shows the residual error per bin, measured in units of the data uncertainty. The shape of the injected signal is also shown.



(a) Fitted vs. injected signal mass at a specified signal strength value.



(b) Fitted vs. injected signal strength at a specified signal mass value.

Figure A11: Fitted values vs. the true values of parameters of the injected signal in the trijet dataset. The bottom panels show the residual error in units of the fitted uncertainty.

Table A6: The candidate functions are obtained from three fits using different random seeds, fitted to the pseudodata of the four-jet spectrum with the (injected) signal region blinded. The fits were performed on a scaled dataset (to enhance fit stability and prevent numerical overflow), and the functions can be transformed back to describe the original spectrum using the transformation: $x \rightarrow 0.000136(x - 1568.5)$. These functions are plotted and compared with the blinded pseudodata in Fig. A13. Numerical values are rounded to three significant figures for display purposes.

	Candidate function (after ROF)	# param.	χ^2/NDF (before ROF)	χ^2/NDF (after ROF)	p-value (after ROF)
SR model 1	$(1.13 \times 10^{-5}x)^{1.28x} / (0.143x^x + 1.63x)$	3	47.95 / 34 = 1.41	39.07 / 34 = 1.149	0.2524
SR model 2	$6.98x^{0.857x} (x + \exp(x))^{-11.8}$	3	47.36 / 34 = 1.393	39.83 / 34 = 1.171	0.2267
SR model 3	$((6.52 \times 10^{-5} + 0.000378x) \tanh(0.641 + 3x))^{x + \tanh(x)} / \tanh(0.145 + x)$	3	71.24 / 34 = 2.095	35.57 / 34 = 1.046	0.3942

Appendix A.4. CMS paired-dijet dataset (1D) [background modeling]

CMS performed a search for high-mass four-jet resonances using proton-proton collision data at a center-of-mass energy of $\sqrt{s} = 13$ TeV and reported no significant deviations from the Standard Model prediction [14]. The dataset for the four-jet spectrum is publicly available on HEPDATA at Ref. [33]. In the analysis, CMS considered four empirical functions to model the background contribution in the distribution of the four-jet invariant mass, m_{jjjj} , and one of them is:

$$f(x) = \frac{p_0(1 - x^{1/3})^{p_1}}{x^{p_2 + p_3 \log x + p_4 \log^2 x}}, \quad (\text{A.3})$$

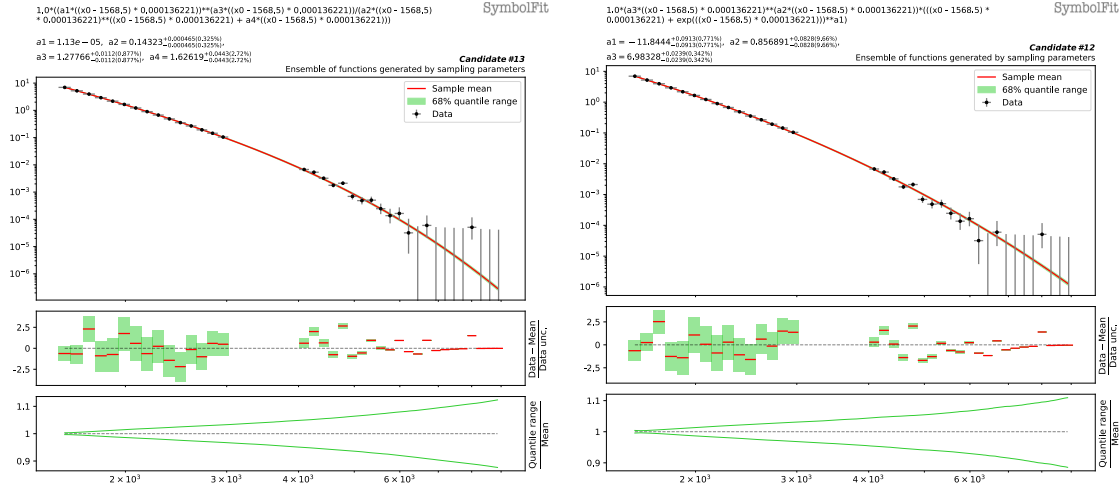
where $x = m_{\text{jjjj}}/\sqrt{s}$ is dimensionless and $p_{\{0,1,2,3,4\}}$ are free parameters.

We perform the same experiments conducted on the dijet dataset, as detailed in Sec. 5.2. Starting from the original four-jet spectrum, we generate pseudodata by injecting a perturbed Gaussian signal centered at $m_{\text{jjjj}} = 3500$ GeV (s_1) with a width of 400 GeV ($2s_2$) and a signal strength of $s_0 = 2$. To model the background, we blind the signal region by masking the m_{jjjj} bins between 3000 and 4000 GeV in the pseudodata and perform the fits.

Three `SymbolFit` runs using different random seeds are carried out, applying the same `PySR` configuration as used for the dijet dataset (see List. 2). Tab. A6 lists the three SR models, each obtained from a fit initialized with a different random seed. The χ^2/NDF scores improve significantly after the ROF step compared to the original functions returned by `PySR`. The three background models fit the blinded pseudodata well, as shown in Fig A12 for the total uncertainty coverage and Fig. A13 for a comparison with the empirical model used by CMS.

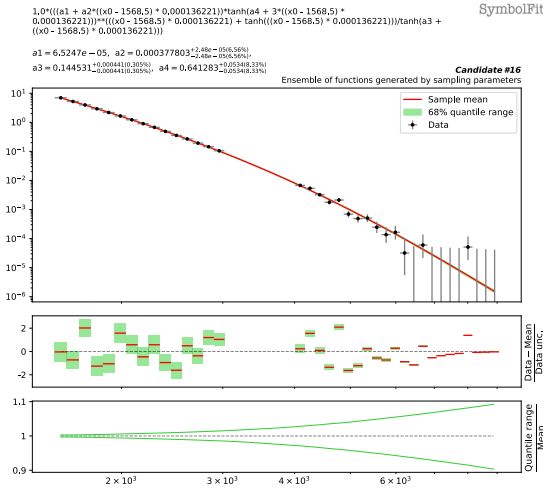
Next, we unblind the pseudodata and perform b-only fits and s+b fits on the full pseudodata spectrum. These results are shown in Fig. A14. In all three SR models, as well as the CMS empirical model, the excess of events over the background around the injected signal location observed in the b-only fits is reduced in the s+b fits, demonstrating that the models are sensitive to the injected signal. Tab. A7 lists the χ^2/NDF scores for each model, showing the fit performance in response to the presence of the injected signal.

To assess whether the SR models can accurately extract the injected signals, we generate multiple sets of pseudodata by injecting Gaussian signals with different mean values ranging from 3350 to 3750 GeV and varying signal strength between 0.5 and 10. We then perform the s+b fits to extract the corresponding signal parameters. Fig. A15 shows the extracted signal parameters plotted against their injected values. All three SR models are capable of extracting the correct signal parameter values within reasonable uncertainties and are comparable to the empirical model used by CMS.



(a) SR model 1.

(b) SR model 2.



(c) SR model 3.

Figure A12: The three SR models fitted to the pseudodata of the paired-dijet spectrum with the signal region blinded (see Tab. A6). To visualize the total uncertainty coverage of each candidate function, the green band in each subfigure represents the 68% quantile range of functions obtained by sampling parameters, taking into account the best-fit values and the covariance matrix within a multidimensional normal distribution. The red line denotes the mean of the function ensemble. At the top of each subfigure, the candidate function and the fitted parameters are shown. The middle panel shows the weighted residual error: $\frac{\text{Data} - \text{Mean}}{\text{Data unc.}}$. The bottom panel shows the ratio of the 68% quantile range to the mean.

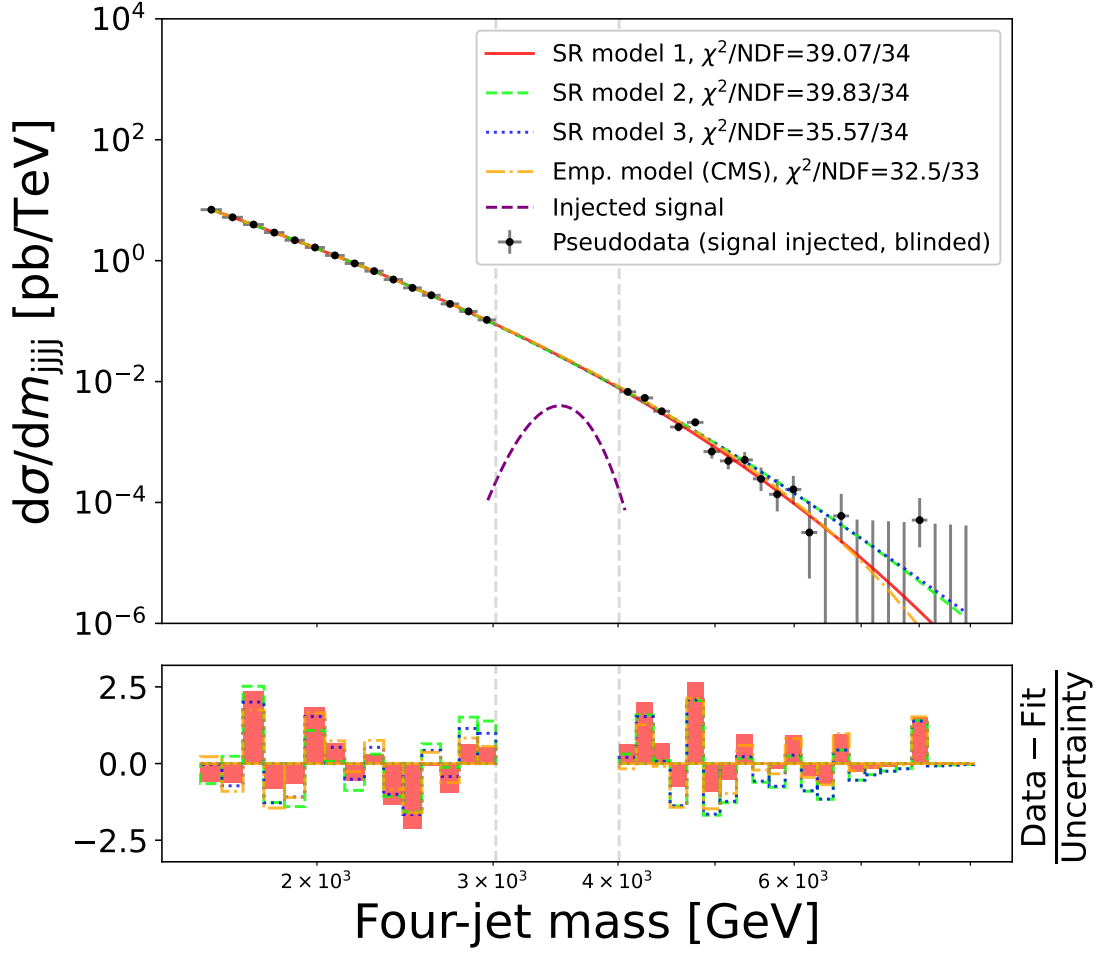


Figure A13: Pseudodata of the four-jet spectrum with the injected signal shown in the blinded signal region. The three SR models (see Tab. A6) are compared against the empirical model used by CMS. The lower panel shows the residual error per bin, measured in units of the data uncertainty.

Table A7: Comparison of the χ^2/NDF scores from three types of fits to the paired-dijet dataset: the b-only fits to the blinded pseudodata, b-only fits to the unblinded pseudodata, and s+b fits to the unblinded pseudodata. The background models used for the fits are listed in Tab. A6, and the fits are shown in Fig. A13 (blinded) and Fig. A14 (unblinded).

	χ^2/NDF (b-only, blinded)	χ^2/NDF (b-only, unblinded)	χ^2/NDF (s+b, unblinded)
SR model 1	39.07 / 34 = 1.149	93.41 / 41 = 2.278	54.38 / 39 = 1.394
SR model 2	39.83 / 34 = 1.171	107.1 / 41 = 2.612	59.14 / 39 = 1.516
SR model 3	35.57 / 34 = 1.046	90.1 / 41 = 2.198	51.88 / 39 = 1.33
Emp. model (CMS)	32.5 / 33 = 0.985	71.84 / 40 = 1.796	44.86 / 38 = 1.181

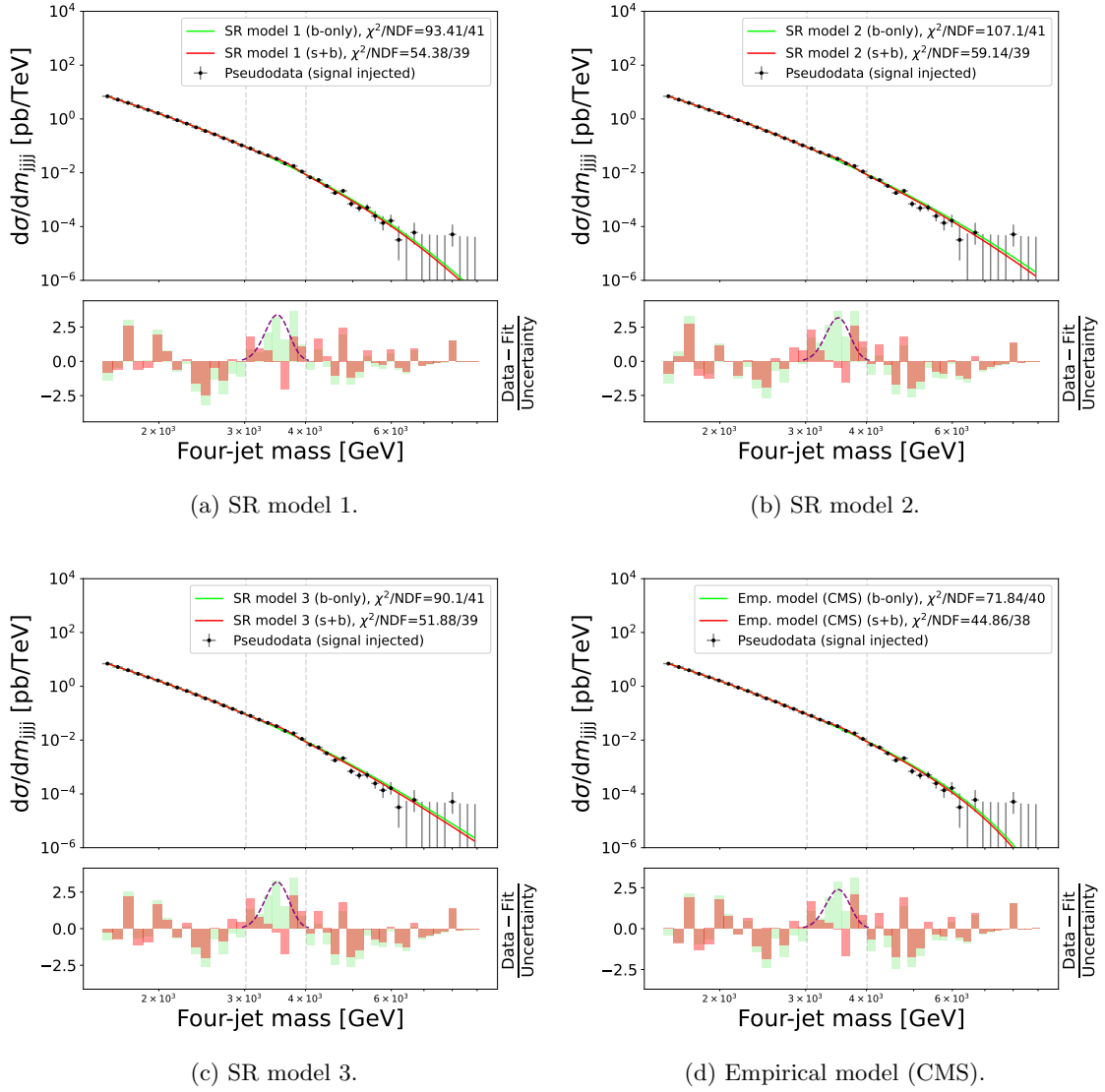
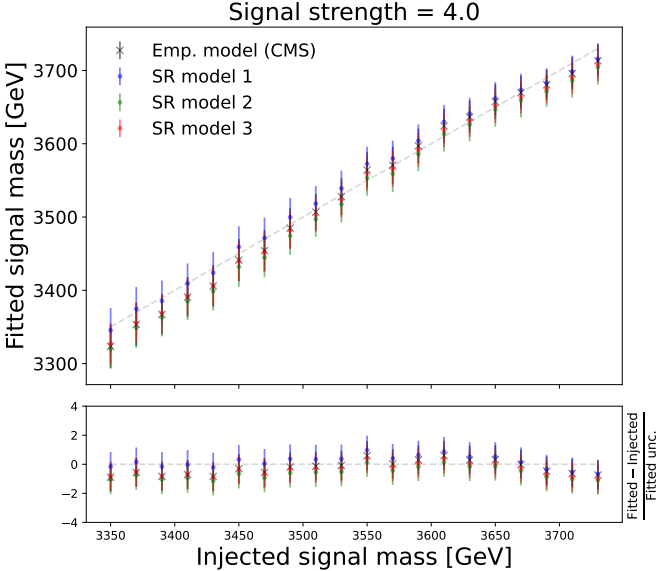
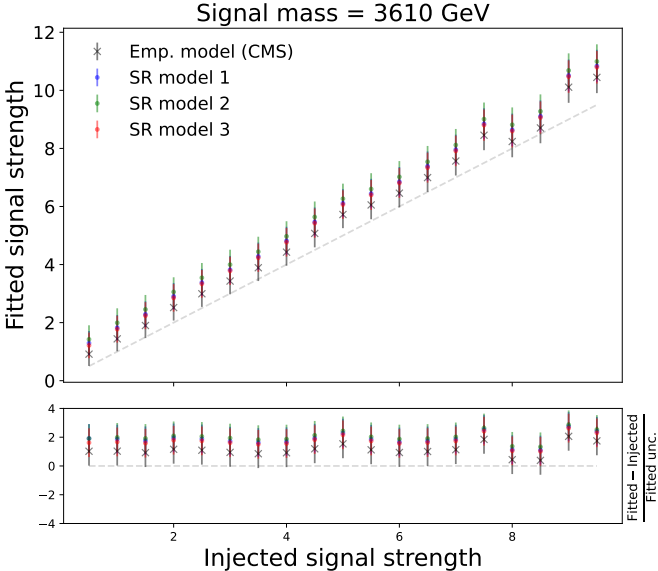


Figure A14: Comparison of the b-only fits and the s+b fits to the unblinded pseudodata of the four-jet spectrum. The lower panel shows the residual error per bin, measured in units of the data uncertainty. The shape of the injected signal is also shown.



(a) Fitted vs. injected signal mass at a specified signal strength value.



(b) Fitted vs. injected signal strength at a specified signal mass value.

Figure A15: Fitted values vs. the true values of parameters of the injected signal in the paired-dijet dataset. The bottom panels show the residual error in units of the fitted uncertainty.

Table A8: The candidate functions are obtained from three fits using different random seeds, fitted to the pseudodata of the dimuon spectrum with the (injected) signal region blinded. The fits were performed on a scaled dataset (to enhance fit stability and prevent numerical overflow), and the functions can be transformed back to describe the original spectrum using the transformation: $x \rightarrow 0.00487(x - 397.4)$. These functions are plotted and compared with the blinded pseudodata in Fig. A17. Numerical values are rounded to three significant figures for display purposes.

	Candidate function (after ROF)	# param.	χ^2/NDF (before ROF)	χ^2/NDF (after ROF)	p-value (after ROF)
SR model 1	$(4.25 \exp(-x))^{1.52+x}$	2	3.469 / 10 = 0.3469	3.007 / 10 = 0.3007	0.9813
SR model 2	$9.68 + x(-7.39 + x)$	2	3.484 / 10 = 0.3484	3.075 / 10 = 0.3075	0.9796
SR model 3	$0.0213x^3 - 5.44 + 3.46$	2	3.456 / 10 = 0.3456	3.066 / 10 = 0.3066	0.9798

Appendix A.5. CMS high-mass dimuon dataset (1D) [background modeling]

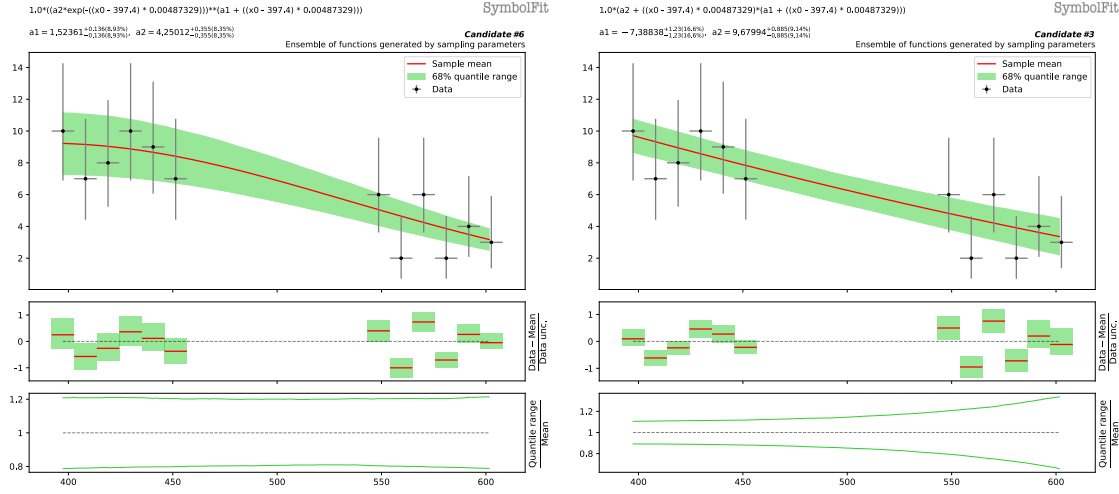
CMS performed a search for high-mass dimuon resonances using proton-proton collision data at a center-of-mass energy of $\sqrt{s} = 13$ TeV and reported no significant deviations from the Standard Model prediction [16]. The dataset for the dimuon spectrum is publicly available on HEPDATA at Ref. [34]. In the analysis, CMS considered three different functions to model the background contribution in the distribution of the dimuon invariant mass, $m_{\mu\mu}$. These functions include a simple exponential, a power-law, and a first-order Bernstein polynomial. Since the dimuon distribution in the signal region is statistically limited, simpler functions are preferred to avoid over-fitting the background. For our comparison, we take the first-order Bernstein polynomial as the empirical model used by CMS.

We perform the same experiments conducted on the dijet dataset, as detailed in Sec. 5.2. Starting from the original dimuon spectrum, we generate pseudodata by injecting a Gaussian signal centered at $m_{\mu\mu} = 500$ GeV (s_1), with a width of 20 GeV ($2s_2$) and a signal strength of $s_0 = 350$. To model the background, we blind the signal region by masking the $m_{\gamma\gamma}$ bins between 450 and 550 GeV in the pseudodata and perform the fits.

Three `SymbolFit` runs using different random seeds are carried out, applying the same `PySR` configuration as used for the dijet dataset (see List. 2), except that the maximum complexity is set at 20 instead of 80, since the diphoton distribution shape is less complex. Tab. A8 lists the three SR models, each obtained from a fit initialized with a different random seed. The χ^2/NDF scores improve significantly after the ROF step compared to the original functions returned by `PySR`. The three background models fit the blinded pseudodata well, as shown in Fig. A16 for the total uncertainty coverage and Fig. A17 for a comparison with the empirical model used by CMS.

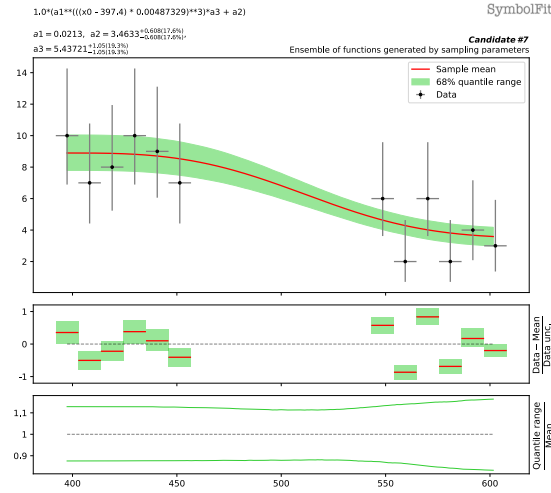
Next, we unblind the pseudodata and perform b-only fits and s+b fits on the full pseudodata spectrum. These results are shown in Fig. A18. In all three SR models, as well as the CMS empirical model, the excess of events over the background around the injected signal location observed in the b-only fits is reduced in the s+b fits, demonstrating that the models are sensitive to the injected signal. Tab. A9 lists the χ^2/NDF scores for each model, showing the fit performance in response to the presence of the injected signal.

To assess whether the SR models can accurately extract the injected signals, we generate multiple sets of pseudodata by injecting Gaussian signals with different mean values ranging from 490 to 510 GeV and varying signal strength between 350 and 600. We then perform the s+b fits to extract the corresponding signal parameters. Fig. A19 shows the extracted signal parameters plotted against their injected values. All three SR models are capable of extracting the correct



(a) SR model 1.

(b) SR model 2.



(c) SR model 3.

Figure A16: The three SR models fitted to the pseudodata of the dimuon spectrum with the signal region blinded (see Tab. A8). To visualize the total uncertainty coverage of each candidate function, the green band in each subfigure represents the 68% quantile range of functions obtained by sampling parameters, taking into account the best-fit values and the covariance matrix within a multidimensional normal distribution. The red line denotes the mean of the function ensemble. At the top of each subfigure, the candidate function and the fitted parameters are shown. The middle panel shows the weighted residual error: $\frac{\text{Data} - \text{Mean}}{\text{Data unc.}}$. The bottom panel shows the ratio of the 68% quantile range to the mean.

signal parameter values within reasonable uncertainties and are comparable to the empirical model used by CMS.

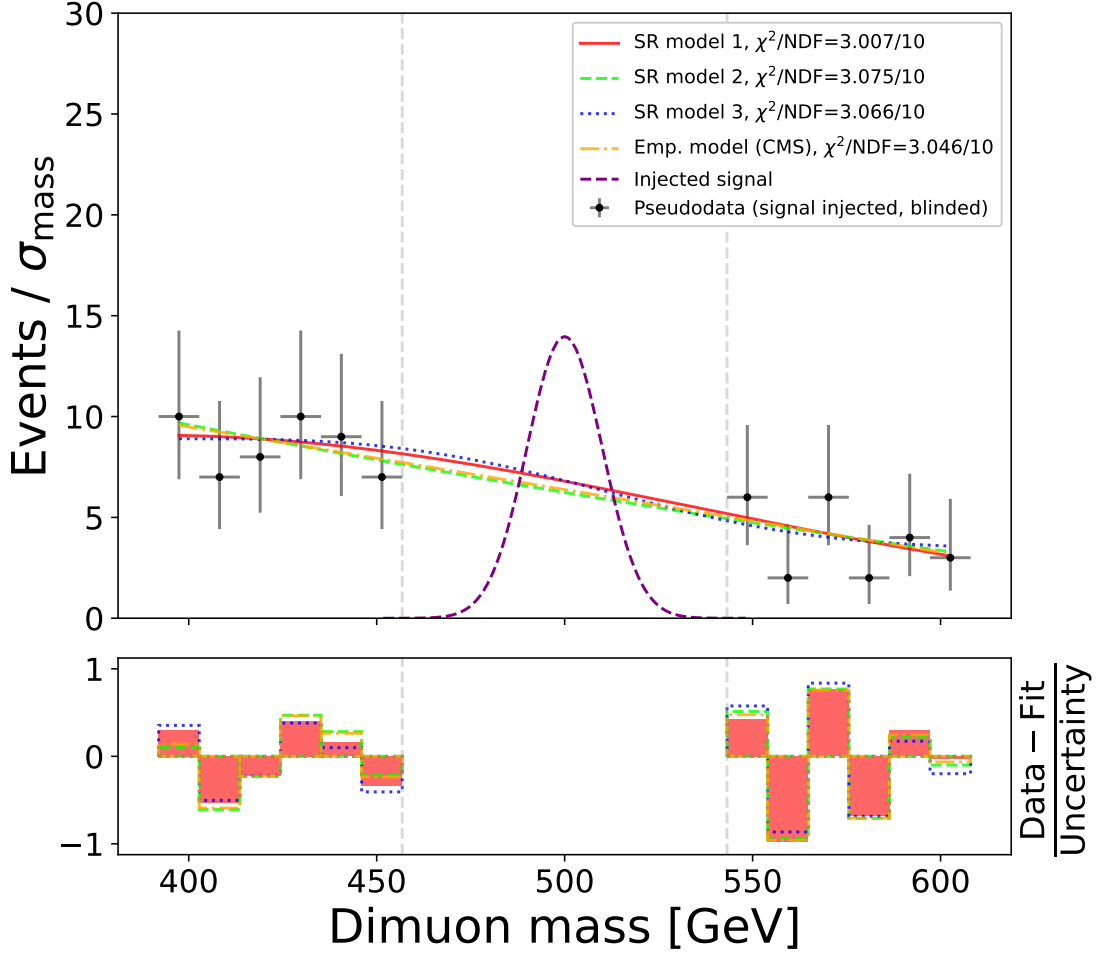


Figure A17: Pseudodata of the dimuon spectrum with the injected signal shown in the blinded signal region. The three SR models (see Tab. A8) are compared against the empirical model used by CMS. The lower panel shows the residual error per bin, measured in units of the data uncertainty.

Table A9: Comparison of the χ^2/NDF scores from three types of fits to the dimuon dataset: the b-only fits to the blinded pseudodata, b-only fits to the unblinded pseudodata, and s+b fits to the unblinded pseudodata. The background models used for the fits are listed in Tab. A8, and the fits are shown in Fig. A17 (blinded) and Fig. A18 (unblinded).

	χ^2/NDF (b-only, blinded)	χ^2/NDF (b-only, unblinded)	χ^2/NDF (s+b, unblinded)
SR model 1	3.007 / 10 = 0.3007	62.98 / 18 = 3.499	15.44 / 16 = 0.965
SR model 2	3.075 / 10 = 0.3075	64.89 / 18 = 3.605	14.32 / 16 = 0.895
SR model 3	3.066 / 10 = 0.3066	62.97 / 18 = 3.498	15.87 / 16 = 0.9919
Emp. model (CMS)	3.046 / 10 = 0.3046	64.43 / 18 = 3.579	14.79 / 16 = 0.9244

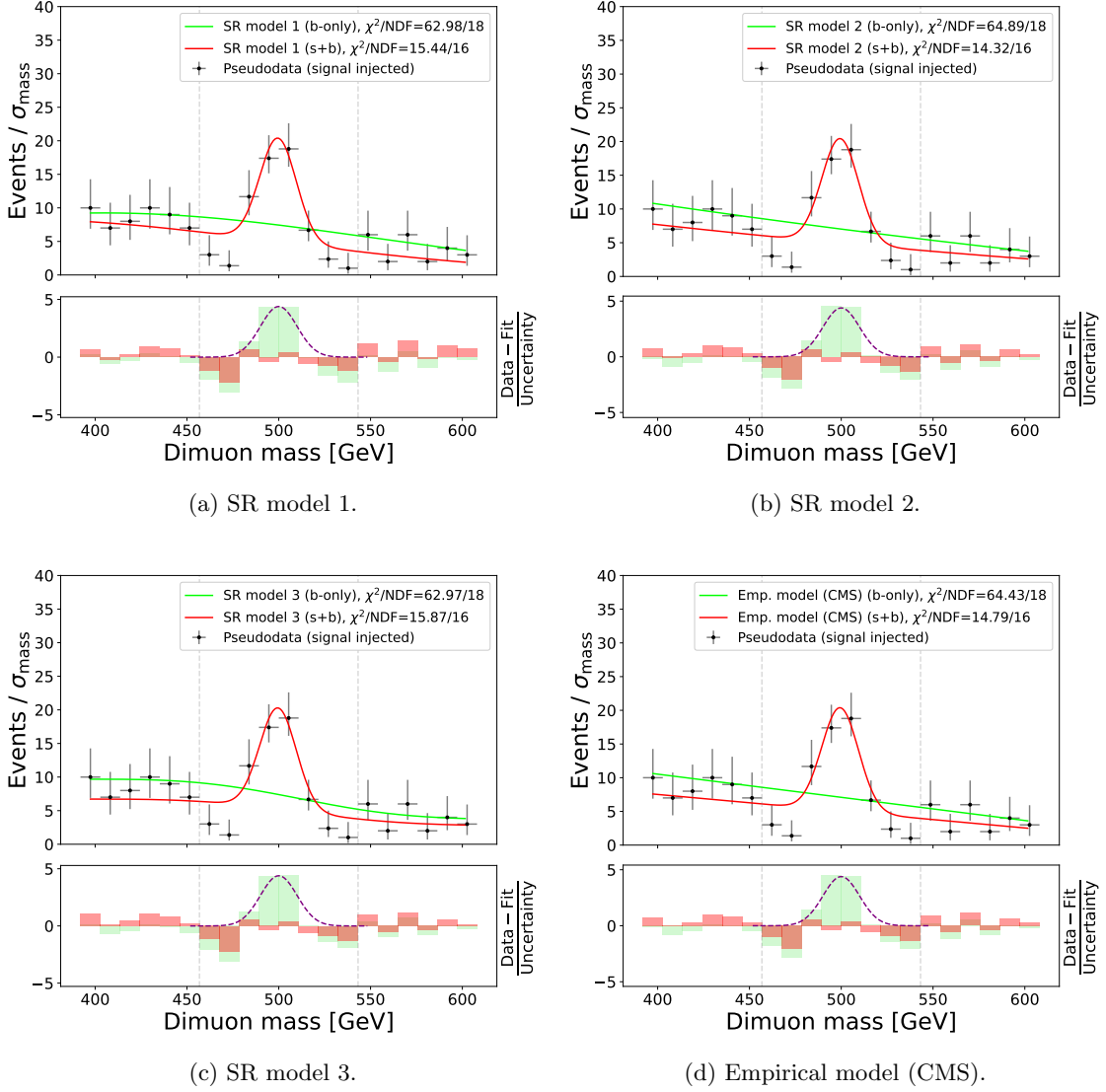
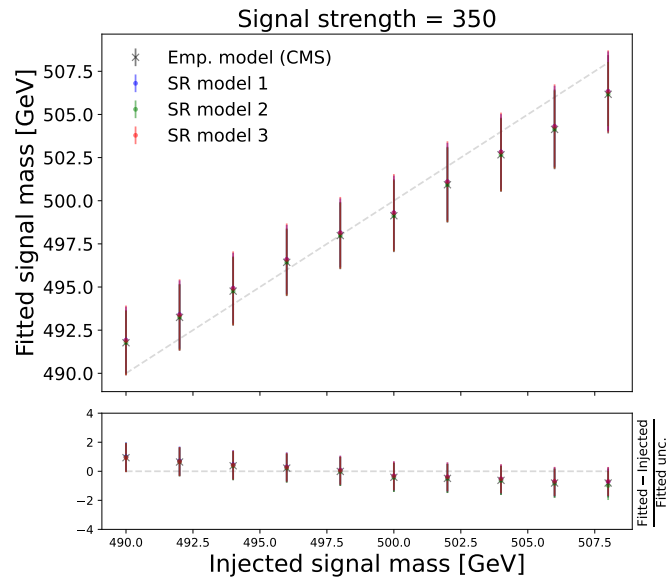
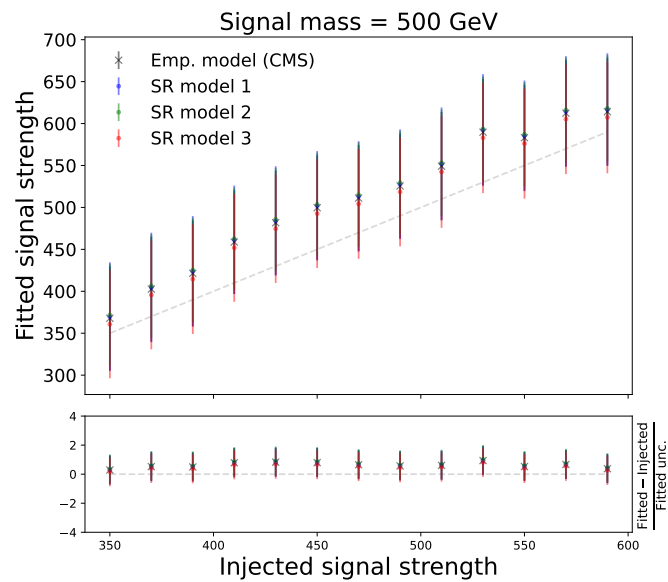


Figure A18: Comparison of the b-only fits and the s+b fits to the unblinded pseudodata of the dimuon spectrum. The lower panel shows the residual error per bin, measured in units of the data uncertainty. The shape of the injected signal is also shown.



(a) Fitted vs. injected signal mass at a specified signal strength.



(b) Fitted vs. injected signal strength at a specified signal mass value.

Figure A19: Fitted values vs. the true values of parameters of the injected signal in the dimuon dataset. The bottom panels show the residual error in units of the fitted uncertainty.