

1 Introduction:

1.1 Motivation:

As there is more data being generated than ever before and new experiments, we need a systematic and automatic way to deduce various mathematical patterns and laws in these data. Through the use of symbolic regression we can utilise these data, and in an explainable manner deduce various new physical laws. In this research I have also extended this beyond physics and have applied this to biological data sets which is a novel application of this method. Perhaps extend this beyond or add a section saying this can also be applied to nlp and that it can learn the rules in language and writing etc. Talk a little about the way this is used outside of this niche use case, and in research, so of course I need to look and research into this.

2 Previos Work:

2.1 Literature Review:

2.2 Introduction:

Humanity has spent millennia observing the world, creating concepts that describe the variables, such as mass and force, to derive laws. In physics, like with all human endeavours, new discoveries and ways of thought are based upon previous works, creating a natural bias in the way new problems are approached. All existing theories, are therefore somewhat biased, this combined with our pre-existing bias in our biological brains, can introduce some hurdles to our future progress [?] [?].

In the 17th Century, Kepler had gotten his hands on the word's most precise data tables on the orbits on planets, using this he spent close to half a decade, and after numerous unsuccessful attempts, he had began a scientific revolution at the time, describing Mar's orbit to be an ellipse [?]. In essence, scientists throughout history, much like Kepler, have spent a great deal of time, discovering the right expressions to match the relevant data they have, this at it's core is symbolic regression. Now, a few centuries later, even with exponential increases in orders of magnitude in our capability to perform calculations through computers, the process of discovering natural laws and the way to express them, has to some extent resisted automation.

One of the core challenges of physics and artificial intelligence, is finding analytical relations automatically, discovering a symbolic expression that accurately matches the data from an unknown function. This problem, due to it's nature, is NP-hard [?] in principle. The vastness of the space of mathematical constants, adds to the difficulty. This literatire review aims to present the recent advances in discovery of emphirical laws through data powered by artificial intelligence. It focuses on methodoogies that diminish human bias through seeking solutions without assumptions. We will explore various techniques employed to achieve these goals, which includes reducing the search space, and analyse the effectiveness of these methods.

Figure 1: This is the orbit of Earth and Mars around the Sun.

2.3 Symbolic Regression:

Symbolic regression, is a technique that analyses and searches over the space of traceable mathematical expressions to find the best fit for a data set. By not requiring prior information about the model, it is unbiased. There are a plethora of various strategies that have been implemented in solving for empirical laws [?], we will explore some of them below. It is also worth mentioning, that unlike other well-known techniques for regression, (eg:

neural networks), that are essentially black boxes, symbolic regression, aims to extract white-box models and is easy to analyse.

Brute Force:

Symbolic Regression (SR), is interpretable [?], unlike Neural Networks (NN), which are often considered more explainable. The difference is interpretability allows us to comprehend how the model works, like observing how gears move in a glass box, while explainable means you get an overview of why a certain output was achieved, even without knowing the full nuances of it's inner workings.

There however, are some challenges associated with SR, in comparison to function fitting (NN). SR, starts with nothing, a blank slate, and it has to learn the entire expression [?], unlike function fitting which just tweaks an already existing function. The exponential search space [?], causes it to be extremely computationally expensive to explore all possibilities. This combined with the fact that, most optimisation algorithms expect a smooth search space [?], however SR lack's smooth interpolation, small changes in the potential solutions (expression), ie: x^3 and $x^3 + 0.1$ can significantly alter the the output. Finally, if the nature of the problem is badly posed [?], there might potentially be multiple solutions to the same data. Imagine trying to find a single polynomial equation with only two points of data, the need to balance finding accurate expressions with finding the most simplistic and generalisable fit, is sometimes troublesome.

The brute force approach of simply trying all possible combinations of symbolic expressions within some defined space. The model will subsequently increase the complexity over time, and will stop when either the fitting errors lowers below some defined limit or exceeds the upper limit of runtime. While in theory can solve all of our problems, in practise takes longer than the age of our universe to finish. In essence it's like searching for a singular drop in the ocean. Thankfully, there are some ways of pruning the search space, and drastically reducing the time taken to solve for the most accurate expression.

Partial Derivatives:

Partial derivatives, of some function f , with multiple variables such as x and y , is it's dervative with respect to one of those two variables, while the other variables in the function are kept constant. Formally, given a function with two or more variables, $f(x_1, x_2, \dots, x_n)$, the partial derivative of f with respect to x_i , where x_i is some value x in $(x_1, x_2, \dots, x_i, \dots, x_n)$, gives the rate of change of f with respect to x_i . It is calculated by taking the i th derivative of f with respect to x_i , whilst holding the other variables fixed. [?] [?]

The partial derivative of a function $f(x, y)$ with respect to x is denoted $\frac{\partial f}{\partial x}$ [?] and is defined:

$$\frac{\partial f}{\partial x} = \lim_{h \rightarrow 0} \left[\frac{f(x+h, y) - f(x, y)}{h} \right]$$

Once you pass in the experimental data, you can pre-process the data, using calculated partial derivatives, for every pair of existing variables. Many physical laws, involve rates of change, and partial derivatives help us represent them. Furthermore it also guides the search process, as the algorithm can use the derivative to accurately represent the underlying laws involved. Through comparing how well the partial derivatives derived through the experimental data compared to the potential expression, the algorithm can assess the accuracy and feasibility of the expressions involved. This strategy can even be extended to prune the search space further, this could be achieved through incorporating knowledge of physics into the constraints for the partial derivatives. These concepts will be illustrated with an example below.

Consider a iron rod, that has been heated up, such that it is hotter on one side than the other. Now it is intuitive to say that closer to the heat source, the temperature will be higher than further along the rod, where it will be colder. We can illustrate this temperature distribution with a function:

$$T(x, y, z)$$

where T is the temperature at a point in the rod, and (x,y,z) are the coordinates along the axis in 3 dimensions. This leads to these 3 partial derivatives:

$$\frac{\partial T}{\partial x}, \frac{\partial T}{\partial y}, \frac{\partial T}{\partial z}$$

These partial derivatives, gives us information about the direction and magnitude of heat flow at various points on the rod. The algorithm then searches for an equation $T(x,y,z)$, that sufficiently predicts the observed temperature distribution and it's partial derivatives, deriving laws such as the heat transfer equations, or elasticity relationships.

$$\frac{\partial T}{\partial t} = \alpha \nabla^2 T$$

Through using partial derivatives, we have in essence redefined the search criteria for the algorithm, through it's measure of the accuracy in comparison of potential solutions over the invariants represented in the experimental data [?]. This also leads to the pleasant finding, that it can additionally capture relationships that represent other identities of the system, beyond invariants and heat transfer equations.

You can subtly guide the type of laws that such an algorithm finds, by selectively picking the variables to input into the algorithm,. For example providing velocities and force to find laws of motion.

Dimensional Analysis:

Dimensional Analysis is a method of solving problems usually in maths and physics, where we analyse the relationships between different physical quantities, by comparing their "units." It is a powerful method of reducing the complexity of systems, enabling engineers and scientists to analyse problems that we can't even pose, much less solve the equations of [?].

Using the fact that numerous questions in science can be simplified by requiring the dimensions/units of the right and left hand side of the expression to be equal, we can transform the question into a smaller number of variables, which all have no dimension [?]. It has been automated to find the integer powers of expressions and has proven to be useful especially when the power is an irrational number.

Here is a general strategy that showcases how dimensional analysis can be used:

Let's say we have a variable in an equation that can be broken down into it's fundamental units, such as (second, kilograms, ampere ...) to various powers. We can then take this, and represent each of the units as vectors, such that each of the fundamental units, is assigned a dimension, and it's important to note, this then allows us to represent any physical quantity as a product of these units, so let us construct a vector v , with 3 integers, where each corresponding integer represents the power of each of the fundamental units.

Given that we want to derive an expression, such as $y = f(x_1, \dots, x_n)$ we can then create some matrix M . Each of the columns of the given matrix, is the unit vector v of the corresponding variable x_i . We then need to define another vector to represent the units of y , which will be called z . If we let the solution be some vector s , solving $Ms = z$, this then lets us raise the powers on both sides, to elevate the independent variables, to make this equation dimensionally consistent.

Taking the null space of the matrix M , where $MV = 0$, allows us a basis to create a dimensionless group, allows for a simplification of the problem.

This is also more intuitive to understand physical phenomena, the nature of physics comprehension, making this vital in further understanding derived laws, making the process easier to explain and understand [?]. Therefore, this is a crucial tool, for cultivating a deeper understanding of physics effectively [?].

Genetic Programming:

Genetic programming (GP), is a special evolutionary algorithmic technique, where the individuals are seen as programs that evolve, starting for a population, is iteratively "evolved," transforming the populations of individual programs, into other populations. This new generation of programs are created using some genetic operations or survival criteria, mimicking natural evolutionary condition on earth.

A very basic overview, shows that genetic programming algorithms, consists of initializing the population, then evaluation of the said population through some predefined metrics and functions, followed by selection of the fittest programs based on the score given by the metric, and "genetic operation," such as reproduction, mutation and cross-over. The algorithm then iterates these steps thousands of times, through many generations, and finally terminates once the desired result has been achieved.

We can use genetic programming, and tweak the algorithm, and combine it with symbolic regression, to help us derive laws.

Consider modelling the various potential formulas as a tree, which is composed of various functions in the nodes. These functions can vary from arithmetic operations , mathematical functions, or defined unique operators. Then we can program the fitness function [?], and use it to measure how well the given potential expression in the population compares with the given databases, and given the nature of genetic programming, the better performing functions are more likely to be passed down into the next generation. Then after many iterations, we can give the solution with the best performance.

There are various ways to implement the fitness function, and for example we can use a criteria like this, along with mean squared error [?]:

$$V = 2X + N \cdot \ln(M/N)$$

Here M is the mean squared error, and N is the number of data points, X is the number of parameters used on the genetic programming algorithm. The lower the value of V is, the better the model performs. The performance of this strategy can then be evaluated with various other metrics, to judge how well the algorithm performs.

3 PySR

This section describes the relevent implementations that are completed as of 10 December 2024.

3.1 Momentum Laws:

To generate the dataset, I chose 100 data points, and created two variables mass (M) and acceleration (a), each represented in two dimensions. Then the data points were generated using *numpy.random.randn* function. The force (F), was then calculated to be the produce of these two data sets. Mass and acceleration were concatenated along the same axis using *numpy.concatenate*, resulting in a combined dataset. This is partially because the model used here, *PySRRegressor* expects a single array as input, and this helps highlight the relationship

between these variables to the symbolic regression algorithm.

Then model performed symbolic regression, configured with 40 iterations along with a customer loss function, taken to be the squared difference between the prediction and the target variable.

$$\mathcal{L}(\hat{x}, x) = (\hat{x} - x)^2$$

The model was trained on this dataset, upon termination, it produced a list of potential candidate formulae, from which I manually identified the correct expression, $F = m\dot{a}$.

Algorithm 1: Symbolic Regression for $F = M \cdot A$

Result: A symbolic representation approximating $F = M \cdot A$

Initialization:

Generate random data for mass (M) and acceleration (A);

Compute target force values: $F = M \cdot A$;

Combine M and A into input matrix X ;

while *Symbolic regression process* **do**

 Train the symbolic regression model with the following settings;;

Binary operators: Multiplication (*);

Unary operators: None;

Loss function: Mean squared error between predictions and targets;

Iterations: 40;

if *Current symbolic representation improves loss* **then**

 Update the symbolic model;

 Save the current best expression;

else

 Continue exploration of new symbolic expressions;

end

end

Similarly, the other laws of momentum, were also dervied using this approach.

$$\mathbf{F}\Delta t = \Delta \mathbf{p} = m(\mathbf{v}_f - \mathbf{v}_i) \quad (1)$$

$$m_1 \mathbf{v}_{1,i} + m_2 \mathbf{v}_{2,i} = m_1 \mathbf{v}_{1,f} + m_2 \mathbf{v}_{2,f} \quad (2)$$

3.2 Pendulum Laws:

The data is generated using numpy. The simulation involves, Euler's method to solve the pendulum's equation of motion. Through taking small and discrete steps, the method approximates the solution. The equation for a simple pendulum is given by:

$$\alpha = -\frac{g}{L} \sin(\theta) \quad (3)$$

α angular acceleration (rad/s²)

g acceleration due to gravity (m/s^2)

L length of the pendulum (m)

θ angular displacement (rad)

The Euler technique approximates the changes in angular velocity and displacement over some small step in time, as follows:

$$\omega_{i+1} = \omega_i + \alpha_i \Delta t \quad (4)$$

$$\theta_{i+1} = \theta_i + \omega_{i+1} \Delta t \quad (5)$$

The function iterates through a few hundred time steps, updating the angular velocity and displacement at each time step. To prevent errors accumulating due to numerical drift, which are small errors that accumulate and become significant due to the inherent nature of approximation methods like Eulers. To keep the values coherent, a wrap around operation is used to ensure the angular displacement is within the range of $[\pi, -\pi]$ radians.

3.3 Noise:

This section investigates the model's robustness to noise. To simulate varying levels of interference, artificial noise was systematically introduced during the data generation phase, building upon the previously established model framework. Noise was modelled by generating random numbers within a range of progressively increasing magnitude, utilizing Python's random library. The objective was to observe and quantify the degradation in model accuracy as a function of increasing noise levels, as well as explore ways to mitigate it.

3.3.1 How noise affects the model:

To introduce noise into the generated dataset, I imported Python's random library and used the randint function. To systematically vary the level of noise, I created an additional function that incrementally increased the parameters passed to randint, causing each successive dataset to become progressively noisier.

After generating these noisy datasets, the symbolic regression model was applied to each one, and the resulting equations were analyzed. The relationship between noise level and model performance was then visualized through a graph. Additionally, the time library was used to measure how long each run of the model took.

PLOT HERE:

3.3.2 Denoise function:

Following the initial noise analysis, a subsequent experiment incorporated a denoising method (implemented using a Python library) applied to the data prior to processing by the pysr model. The results, as depicted in the accompanying plot, indicate that this denoising approach improved model performance up to a specific noise threshold. However, beyond this critical level, performance degraded comparably for both the denoised and non-denoised datasets, suggesting the denoise function's effectiveness diminished at higher noise intensities.

PLOT HERE:

4 Symbolic Regression from First Principles:

4.1 A Brute-Force Approach:

The core and essential component of any symbolic regression model lies in its ability to generate and traverse the search space of potential equations and expressions that best fit the given data. To streamline the process and validate the functionality of my expression generation, I began by implementing simple two-variable equations, specifying the operations used within the equation. This approach was initially limited to basic operations, with plans to extend it to accommodate constants and additional complexities.

Algorithm 2: Generate Initial Symbolic Expressions

Data: List of variable names, list of constants, list of operators

Result: List of generated symbolic expressions

Initialization::

Create pool of symbolic variables and constants from inputs;

Initialize empty list for expressions;

```
for each pair (a, b) from the pool (considering permutations) do
    for each operator op from the list of operators do
        Form expression a op b;
        if resulting expression is valid then
            Add expression to the list of expressions;
        end
    end
end
```

end

Return list of expressions;

Subsequently, I refined this approach by designing a recursive method to generate expressions, allowing for the creation of more robust and diverse equations from the available variables. This process is dynamic, making it adaptable to a variety of input configurations.

Algorithm 3: Recursive Generation of Symbolic Expressions

Data: List of operators, list of variable names, maximum depth

Result: List of symbolic expressions generated at the maximum depth

Initialization::

Create list 'Expressions[0]' containing symbolic variables from input;

```
for depth d from 1 to maximum depth do
    Create empty list 'Expressions[d]';
    for each expression a in 'Expressions[d-1]' do
        for each expression b in 'Expressions[d-1]' do
            for each operator op in list of operators do
                Attempt to form new expression a op b;
                if resulting expression is valid then
                    Add new expression to 'Expressions[d]';
                end
            end
        end
    end
end
```

end

Return 'Expressions[maximum depth]';

4.1.1 Exploiting Physical Properties:

The next step involves truncating the generated expressions to prune the search tree as efficiently as possible. One effective method for achieving this is by leveraging the symmetrical properties of physical equations and recognizing their mathematical equivalence. This includes removing redundant or duplicate expressions that do not contribute new information.

This is the approach I used to achieve this:

4.1.2 Dealing with constants:

Another approach I employed to further prune the set of generated expressions was by filtering out any expressions that did not contain all the specified variables. This step helps optimize the process by reducing the number of irrelevant expressions, ultimately saving computation time during the evaluation phase.

4.1.3 Dealing with powers:

To apply powers to expressions, I incorporated the power operation directly into the generated expressions. This not only allows for the creation of more complex models but also facilitates further pruning of the search tree by avoiding the generation of redundant expressions that already incorporate powers.

Algorithm 4: Apply Powers Recursively to Expressions

Data: List of initial expressions, list of powers, maximum depth

Result: List of unique symbolic expressions generated by applying powers

Initialization:

Create set 'ResultExpressions' and add all initial expressions to it;

for depth d from 1 to maximum depth **do**

 Create empty set 'NewExpressionsThisDepth';

for each expression e in 'ResultExpressions' (from previous depths) **do**

for each power p in list of powers **do**

 Attempt to calculate $\text{powered_}e = e^p$;

if resulting expression is valid **then**

 Add $\text{powered_}e$ to 'NewExpressionsThisDepth';

end

end

end

 Add all expressions from 'NewExpressionsThisDepth' to 'ResultExpressions';

end

Return list representation of 'ResultExpressions';

I implemented a filtering mechanism based on whether the expression included a power operation. This further reduces the size of the search tree by eliminating unnecessary branches. While a more robust model would derive this power operation from scratch, I opted for this approach to optimize computation time and maintain flexibility.

Algorithm 5: Filter Expressions by Presence of Target Powers

Data: List of expressions, list of target powers

Result: List of expressions containing a sub-expression with a target power

Initialization:

Create empty list 'FilteredExpressions';

for each expression e in input expressions **do**

 Set 'found_{target_ppower}' to false *Traverse the symbolic tree of e (e.g., preorder);*

for each sub-expression s during traversal **do**

if s is a power AND exponent of s is in target powers **then**

 Set 'found_{target_ppower}' to true *Break traversal for e ;*

end

end

if 'found_{target_ppower}' is true **then**

 Add e to 'FilteredExpressions';

end

end

Return 'FilteredExpressions';

4.1.4 Chaining powers and constants:

The next step was to chain together powers and constants, applying both to the generated expressions. While the existing design already supports chaining, it is essential to properly filter the results to ensure the search tree remains as compact as possible.

Although constants can be filtered using the existing method, the power operations are embedded within the constants, which causes the previous power filter to no longer function as expected. Consequently, I redesigned the filtering process to operate recursively, allowing it to handle both powers and constants effectively.

code:

However, this approach sometimes results in expressions that feature chained constants, such as $\sin(\sin())$. To refine the model further, I introduced an additional filter to remove expressions with multiple instances of the same constant chained together.

4.1.5 Loading data:

Next, I needed an efficient way to load the data I had created. At this stage, my primary focus was on rapid testing. To facilitate this, I initially generated some dummy data values. Afterward, I decided to store the data as a NumPy array, as this would offer significant speed advantages over using text files. Several factors contribute to this, such as NumPy arrays being stored in memory, the efficiency of the underlying binary data format, and NumPy's use of C, which allows it to vectorize operations, greatly enhancing performance.

Algorithm 6: Load and Validate Data

Data: Input matrix X , target vector Y , list of variable names

Result: Input matrix X , target vector Y , list of symbolic variables

Check Input Validity::

if *number of columns in X \neq number of variable names* **then**

 | Indicate error or stop execution;

end

Process Variables::

Create list of symbolic variables from the list of variable names;

Return X , Y , and the list of symbolic variables;

As shown, I perform a check to ensure that the number of variables provided matches the shape of the array X , where X represents the input data, and y represents the target data, i.e., the final result. For example, X contains the mass and acceleration values, while y contains the corresponding values of the force f as calculated by the equation $f=ma$. This serves as a basic validation to confirm that the number of columns in the input data corresponds correctly to the variables provided.

4.1.6 How to mitigate noise in data:

Ways to mitigate the noise and its affects on the model were explored. Functions such as "denoise," in the symbolic regression library helped to some extent. However after a certain point, such methods do not seem to offer much assistance.

I also made my own denoise algorithm. I implemented various different denoise algorithms to see what effects they had. Firstly I implemented a simple moving average as a way to mitigate the noise in the dataset. reword this -> "Simple and fast, smooths data well by averaging neighbors. However, it blurs sharp changes and is sensitive to extreme outlier values, pulling the average significantly and distorting the signal." These were my results, this is the pseudo code, explain the algorithm

The second denoise algorithm I implemented is a median filter, and this is what effects it has, and this is how i implemented it. Insert Pseudo code. reword: "Excellent at removing spikes and preserving edges better than averaging. Less affected by outliers. Can sometimes slightly distort the overall shape of the signal, especially with large window sizes."

Finally this is the third algorithm that I had implemented for denoising. Wavelet Denoising, this is the effects, and this is the pseudo code. Reword this -> "Transforms data to isolate noise, preserving both smooth and sharp signal features effectively. More complex to understand and requires careful selection of wavelet type and parameters for optimal results, which can be tricky."

4.1.7 Evaluating expressions:

Next, I evaluate the expressions that have been generated. I assign the input variables to the corresponding columns of the data in increasing order. These values are then substituted into the expressions, and the model runs the calculations, producing an array of outputs for each expression. This process essentially evaluates every pruned expression and returns a NumPy array of results based on the input data.

Algorithm 7: Evaluate Symbolic Expression Numerically

Data: Symbolic expression, list of variable names, input data matrix X

Result: Numerical evaluation of the expression for each data point in X , or NaN if evaluation fails

Initialization::

Create list of symbolic variables from the list of variable names;

Convert to Numerical Function::

Convert the symbolic expression into a numerical function, mapping symbolic variables to input columns of X ;

Evaluate::

Prepare input data from columns of X to match function arguments;

Evaluate the numerical function using the prepared input data;

Return the resulting array of values;

Evaluation fails (e.g., division by zero) Return an array of Not-a-Number (NaN) values with the same size as the number of data points;

Following the approach outlined in the paper [insert citation here], I utilized a medium error description length loss function, implementing it as described. The error is calculated using the squared difference to ensure all errors are positive, and a constant of 1 is added to guarantee that all errors are greater than 1 when taking the logarithm.

Algorithm 8: Calculate Mean Log Squared Error

Data: Array of predicted values, Array of original values

Result: Mean of the base-2 logarithm of (1 + squared error)

Initialization::

Flatten input arrays ‘predicted’ and ‘original’;

Initialize ‘total_log_error’ to 0.0 *Get the number of elements ‘n’ (length of arrays)*

Calculate Total Log Error::

for index i from 0 to $n - 1$ **do**

 Calculate ‘error’ = absolute difference between ‘original[i]’ and ‘predicted[i]’;

 Calculate ‘squared_error’ = ‘error’ squared Calculate ‘log_error’ =

 base - 2 logarithm of (1 + ‘squared_error’) Add ‘log_error’ to ‘total_log_error’ **end**

Calculate Mean Error::

 Calculate ‘mean_error’ = ‘total_log_error’ / ‘n’

 Return ‘mean_error’

4.2 Polynomial Fit Module:

Now that the core of the algorithm is functional—handling constants, powers, variables, generating expressions, and filtering redundancy using physical principles like symmetry—I aimed to extend the program by implementing a polynomial fitting module. The goal of this technique is to efficiently fit data to a polynomial model, as many functions in physics (or parts of them) are well-approximated by low-order polynomials, and polynomial fitting is a computationally inexpensive method for this specific class of functions. The technique generates all possible polynomial terms up to a specified degree (e.g., degree 4) and creates a linear equation for each data point where the unknowns are the polynomial coefficients. The system of equations is solved using standard methods such as least squares, and the Root Mean Squared Error (RMSE) of the fit is calculated. If the RMSE is below a predefined tolerance (denoted as pol), the polynomial is accepted as a solution. This approach serves as a fast base case in the recursive algorithm, quickly solving problems that are simple polynomials, and it can also handle sub-problems transformed into polynomial form by other modules, such as dimensional analysis or function inversion.

4.2.1 Data Loading:

To begin, I developed the data loading function. The goal was to accept a NumPy array containing the data, along with a list of variables. The function then compares the shape of the data array with the number of variables provided to ensure that the input is consistent and sufficient for further processing.

Algorithm 9: Load and Validate Data Array

Data: Input data array, list of variable names

Result: Validated input data array (if valid)

Check Input Validity::

Get number of columns in data array;

Get number of variable names;

if *number of columns* \neq *number of variable names* **then**

 Raise a value error indicating the mismatch;

end

Return input data array;

4.2.2 Generating polynomial expressions:

The next step involves generating polynomial expressions. The function returns a list of polynomial expressions based on the input coefficients, variables, and operators, considering a specified maximum degree for the terms.

The function works by first creating symbolic representations for the variables. It then iterates over all possible combinations of powers for the variables up to the specified degree and combines these terms using the provided operators. Finally, the generated expressions are simplified and returned as a list.

Algorithm 10: Generate Expressions from Terms and Operators

Data: List of coefficients, list of variable names, list of operators, maximum degree

Result: List of generated symbolic expressions

Initialization::

Create list of symbolic variables from variable names;

Initialize empty list 'GeneratedExpressions';

Generate Expressions::

for *each combination of powers (from 1 to max degree) for each variable* **do**

 Create a list of 'terms', where each term is 'coefficient * variable

^{Power} **for** *each combination of operators (one less than number of terms)* **do**

 Build an 'expression' by combining the 'terms' sequentially using the chosen operators;

 Simplify the 'expression';

 Add the simplified 'expression' to 'GeneratedExpressions';

end

end

Return 'GeneratedExpressions';

4.2.3 Filtering the Polynomial expressions:

The `filter_expressions` function programmatically filters symbolic expressions based on both structural and semantic constraints. It scales symbolic filtering tasks where strict mathematical structures must be enforced.

The initial version of this function worked for symbolic constants (e.g., sin, cos, etc.) but failed to handle numbers or integer coefficients. This issue was identified during testing, prompting me to rewrite the function so

that it could also handle integer coefficients properly.

Algorithm 11: Filter Expressions by Variables, Constants, and Power

Data: List of expressions, List of required variable names, List of required constants (values/types),
Required power value

Result: List of expressions matching all criteria

Initialization::

Create empty list 'FilteredExpressions';

Convert required variable names to symbolic variables;

for each expression e in input expressions **do**

 Set ' $vars_{ok}$ ' = $true$ if all required symbolic variables are in e 's free symbols, false otherwise;

if NOT ' $vars_{ok}$ ' **then**

 | continue

end

 ;

 Set ' $constants_{ok}$ ' = $true$ **for** each required constant c in list of required constants **do**

 Check if c is present as a sub-expression in e (matching value or type);

if NOT present **then**

 | ' $constants_{ok}$ ' = $false$; break loop over constants **end**

 ;

end

if NOT ' $constants_{ok}$ ' **then**

 | continue

end

 ;

 Set ' $power_{ok}$ ' = $false$ Check if any sub-expression in e is a power with exponent equal to required power;

if present **then**

 | ' $power_{ok}$ ' = $true$ **end**

 ;

if NOT ' $power_{ok}$ ' **then**

 | continue

end

 ;

 Add e to 'FilteredExpressions';

end

Return 'FilteredExpressions';

4.2.4 Evaluating expressions:

The next step involves fitting the filtered expressions to the dataset. The model fitting function fits polynomial expressions to the input data by determining the optimal set of coefficients that minimize the error between the predicted and actual output values. It evaluates multiple polynomial degrees, up to a specified maximum, and selects the degree that results in the lowest error, thereby ensuring an optimal balance between accuracy and complexity.

I utilize the Root Mean Squared Error (RMSE) to calculate the loss, and the function returns a list of loss values, one for each expression.

Algorithm 12: Evaluate Expressions and Calculate RMSE

Data: List of symbolic expressions, List of symbolic variables, Input data matrix X , True target vector Y_{true}

Result: List of (expression, RMSE) pairs for successfully evaluated expressions

Initialization:

Create empty list 'EvaluationResults';

for each expression e in input expressions **do**

 Convert e into a numerical function using the symbolic variables;

 Evaluate the numerical function for each data point in X to get predicted values Y_{pred} ;

 Calculate Root Mean Squared Error (RMSE) between Y_{pred} and Y_{true} ;

 Add the pair (e , RMSE) to 'EvaluationResults';

 Evaluation fails (e.g., runtime error)

end

Return 'EvaluationResults';

To begin the fitting process, I take an expression, substitute the input variables with the corresponding values from the dataset, and compute the predicted yy -values based on the equation. I then calculate the difference between the predicted values and the true target values (yy), which are the actual outputs. This difference is used to compute the Root Mean Squared Error (RMSE), which quantifies the prediction error for the expression.

4.2.5 Best Polynomial Fit:

After calculating the RMSE values for each expression, I select the expression with the lowest RMSE as the most accurate polynomial fit for the data. This ensures that the chosen model has the best performance in terms of minimizing prediction error.

```
def bestFit(results): return min(results, key=lambda x: x[1])
```

4.3 Dimensional Analysis:

Physical equations must be dimensionally consistent, meaning the units on both sides of the equation must match, which severely limits the possible forms of the unknown function. This dimensional constraint provides a strong simplification of the problem, significantly narrowing the scope of valid equations. AI Feynman addresses this by applying dimensional analysis as the first step, simplifying the problem by identifying which combinations of variables are dimensionally consistent. The units of the variables—such as mass, length, and time—are represented as vectors of integer powers, forming a linear system based on the unit vectors of the input and target variables. Solving this system and finding the null space reveals dimensionless combinations of variables, which transforms the problem into one of finding a function of these new dimensionless variables. This process typically reduces the number of independent variables that the algorithm needs to search over, drastically shrinking the combinatorial search space for subsequent steps, such as polynomial fitting, brute force, and neural network-guided searches. As a result, the reduction in variables leads to a significant boost in efficiency, making these searches faster and more likely to succeed.

4.3.1 Handling Units:

The AI Feynman database was accessed, and the units.csv file was downloaded to better understand the units present in the dataset. Upon reviewing the required units, a unit table was created in the form of an array, where each unit corresponds to a unique power of the fundamental SI units. Additionally, the basic SI units were implemented as an array/list to facilitate this mapping.

Algorithm 13: Dimensional Vectors for SI Base Units

Mapping from physical quantity name to its dimensional vector (Mass, Length, Time, Temperature, Current, Amount, Luminous Intensity)::

mass: [1, 0, 0, 0, 0, 0, 0];

length: [0, 1, 0, 0, 0, 0, 0];

time: [0, 0, 1, 0, 0, 0, 0];

temperature: [0, 0, 0, 1, 0, 0, 0];

current: [0, 0, 0, 0, 1, 0, 0];

amount: [0, 0, 0, 0, 0, 1, 0];

luminous_intensity: [0, 0, 0, 0, 0, 0, 1];

There were also relevant derived units included.

4.3.2 Construct Matrix and Target Vector:

This function constructs the dimensional matrix MM and the target vector bb , which are essential for performing dimensional analysis. It accepts lists of independent and dependent variable names, along with a dictionary that maps each variable name (as the key) to its corresponding unit vector (as the value). The unit vectors for the independent variables are retrieved through dictionary lookup, using lowercase variable names (i.e., `var.lower()`) to ensure case-insensitivity. These vectors are then efficiently assembled into the columns of matrix MM using `numpy.column_stack`, while the unit vector of the dependent variable forms the target vector bb . This approach ensures usability, case-insensitivity, leverages the performance benefit of `numpy.column_stack`, and includes explicit error handling to prevent issues.

Algorithm 14: Construct Dimensional Matrix and Target Vector

Data: List of independent variable names, Dependent variable name, Dictionary mapping variable names to dimensional vectors

Result: Matrix M of independent variable dimensional vectors, Vector b of dependent variable dimensional vector

Construct Matrix M :

Initialize empty matrix M ;

for each independent variable name v in the list **do**

 Look up the dimensional vector for v in the dictionary;

 Add this vector as a column to matrix M ;

 Variable name v not found in dictionary Raise an error indicating the missing independent variable;

end

Construct Vector b :

Look up the dimensional vector for the dependent variable name d in the dictionary;

Set this vector as vector b ;

Variable name d not found in dictionary Raise an error indicating the missing dependent variable;

Return matrix M and vector b ;

4.3.3 Solving Dimension and Basis Units:

The `solveDimension` function solves the system of equations $Mp=b$ for the unknown vector pp , where MM is the dimensional matrix and bb is the target vector. The function begins by converting the input matrices MM and bb into symbolic matrices using SymPy's `Matrix` class. It then attempts to solve for pp using the LU decomposition method (`LUsolve`), which is efficient for solving linear systems. If this process fails, an error is raised, indicating the issue encountered during the solution attempt. Additionally, the function computes the null space of matrix MM , representing the set of dimensionless combinations of the variables. The function returns two outputs: the solution vector pp and the null space UU , which provides insight into any dimension-

less combinations of the input variables. This method ensures robust error handling and leverages symbolic computation for accuracy.

Algorithm 15: Solve Dimensional System and Find Null Space

Data: Matrix M , Vector b

Result: Particular solution vector p , Null space basis matrix U

Solve for Particular Solution p :

Solve the linear system $M \cdot p = b$ for vector p ;

Solving fails (e.g., matrix is singular or system inconsistent) Raise an error indicating failure to solve the system;

Calculate Null Space Basis U :

Calculate the basis vectors for the null space of matrix M ;

Form matrix U where columns are the null space basis vectors;

Return p and U ;

4.3.4 Data Transformation Function:

The generate_{dimensionless}data function is designed to transform a dataset into dimensionless form by applying the scaling factors.

If the null space U is provided, the function proceeds to generate new dimensionless variables by computing the product of the input data $data_{x \times p}$ raised to the powers specified by each vector in U . These newly generated dimensionless variables are stacked together to form a transformed input dataset $data_{x \times p}^{prime}$. If no null space is provided, the original input data is used.

Algorithm 16: Generate Dimensionless Data

Data: Input data matrix X , Target data vector Y , Particular solution vector p , Null space basis matrix U

Result: Dimensionless input matrix X' , Dimensionless target vector Y'

Prepare Particular Solution:

Ensure p is a flattened numerical vector;

Calculate Scaling Factor:

Calculate a 'scaling_factor' vector by raising each column of X to the corresponding power in p and taking the product across variables for each data point;

$scaling_factor_i = \prod_j X_{ij}^{p_j}$;

Transform Target Data:

Calculate dimensionless target vector Y' by dividing Y by the 'scaling_factor' $Y'_i = Y_i / scaling_factor_i$;

Transform Input Data:

if Null space basis U is not empty **then**

 Create empty list 'DimensionlessVariables';

for each vector u in U **do**

 Calculate a new dimensionless variable vector by raising each column of X to the corresponding power in u and taking the product across variables for each data point;

$new_var_i = \prod_j X_{ij}^{u_j}$;

 Add 'new_var' to 'DimensionlessVariables' **end**

 Form dimensionless input matrix X' by stacking the vectors in 'DimensionlessVariables';

end

else

 Set dimensionless input matrix X' equal to the original input matrix X ;

end

 Return X' and Y' ;

4.3.5 Symbolic Transformation Generator:

This function generates the symbolic mathematical expressions corresponding to the dimensional analysis transformation. It accepts the original independent variable names (*independent_vars*), *the exponent vectors for scaling(p)*, and *the dimensionless combinations*.

The function first creates symbolic representations of each independent variable using `sp.symbols` from the SymPy library. Then, using the scaling exponents `pp`, it constructs the symbolic expression *symbolic_p representing the unit-fixing scaling factor $x p x$ through $sp.Mul$* , which allows for the multiplication of terms. This expression effectively represents the scaling factor.

Next, the function iterates through each exponent vector `uu` in the null space `UU`, building the corresponding symbolic expressions for the dimensionless combinations. Each new variable is constructed by applying the powers from the exponent vector `uu` to the original input variables, and the resulting expressions are added to a list. This process ensures that both the scaling factors and the dimensionless combinations are represented as symbolic mathematical expressions, which are essential for understanding the relationship between the variables in the dimensional analysis.

Algorithm 17: Symbolic Transformation using Dimensional Analysis Results

Data: List of independent variable names, Particular solution vector p , Null space basis matrix U

Result: Symbolic scaling factor expression, List of symbolic dimensionless group expressions

Initialization::

Create list of symbolic variables from independent variable names;

Construct Symbolic Scaling Factor::

Form symbolic expression for scaling factor by taking the product of each symbolic variable raised to the corresponding power in p ;

$\text{ScalingFactor} = \prod_i \text{variable}_i^{p_i}$;

Construct Symbolic Dimensionless Groups::

Initialize empty list 'SymbolicDimensionlessGroups';

for each vector u in Null space basis U **do**

 Form symbolic expression for a dimensionless group by taking the product of each symbolic variable raised to the corresponding power in u ;

$\text{DimensionlessGroup}_u = \prod_i \text{variable}_i^{u_i}$;

 Add 'DimensionlessGroup _{u} ' to 'SymbolicDimensionlessGroups' **end**

 Return 'ScalingFactor' and 'SymbolicDimensionlessGroups';

4.4 Neural Network Fitting:

The next critical component involves using neural networks to predict and compute gradients from the dimensionless data, providing valuable insights for visualization. While neural networks do not directly solve for symbolic expressions, they serve as powerful tools for approximating complex relationships within the data. By training a neural network on the dimensionless variables, it can predict the output for a given input and calculate gradients, offering a smooth, differentiable function that helps visualize how changes in the input variables influence the model's predictions. This capability allows us to gain a deeper understanding of the underlying behavior of the system. Although neural networks do not offer an explicit symbolic expression, they provide a flexible and efficient way to visualize the functional dependencies, aiding in the interpretation of complex patterns that may be difficult to express symbolically. Ultimately, the network's predictions and gradients can be used to explore and understand the data, even if the true functional form remains implicit.

4.4.1 SymbolicNetwork:

This class defines the neural network architecture used as a universal function approximator within the symbolic regression framework. It inherits from `torch.nn.Module`, the base class for all neural network modules in Py-

Torch. The constructor ($init$) initializes the network structure, accepting the number of input features (n_{input}) and defaulting to a single output (n_{output}). The model

Algorithm 18: Symbolic Regression Neural Network Architecture

Class Definition::

Define a Neural Network Class ‘SymbolicNetwork’ inheriting from a base Neural Network Module;

Initialization Method ($init$): *Number of input features* (n_{input}), *Number of output features* (n_{output} , default 1)

;

Call the constructor of the base Neural Network Module;

Define a sequential model containing the following layers::

Linear layer: n_{input} inputs, 128 outputs Tanh activation function Linear layer :

128 inputs, 128 outputs Tanh activation function Linear layer :

128 inputs, 64 outputs Tanh activation function Linear layer :

64 inputs, 64 outputs Tanh activation function Linear layer : 64 inputs, n_{output} outputs

Forward Method (**forward**):;

Data: Input tensor ‘x’

;

Result: Output tensor

;

Pass the input tensor ‘x’ through the defined sequential model;

Return the output tensor;

4.4.2 Preparing the data:

This function preprocesses raw input ($data_x$) and output ($data_y$) data into a format suitable for PyTorch model training and validation. Splitting is crucial for monitoring generalization and model performance.

Algorithm 19: Prepare Data for Neural Network Training

Data: Input data array X , Target data array Y , Batch size, Training split ratio (default 0.8)

Result: Training data loader, Validation data loader

Convert to Tensors::

Convert X to a PyTorch tensor (float32);

Convert Y to a PyTorch tensor (float32) and add a dimension;

Split Data::

Get the total number of samples;

Calculate the index for splitting based on the training split ratio;

Split the input and target tensors into training and validation sets;

Create Datasets::

Create a training dataset from the training tensors;

Create a validation dataset from the validation tensors;

Create Data Loaders::

Create a training data loader from the training dataset, using the specified batch size and enabling shuffling;

Create a validation data loader from the validation dataset, using the specified batch size and disabling shuffling;

Return the training data loader and the validation data loader;

4.4.3 Training the Network:

This function orchestrates the supervised training process for the provided PyTorch neural network model, with the primary goal of adjusting the model’s parameters (weights and biases) to minimize the difference between its predictions and the true target values using the training data, while also monitoring performance on unseen validation data. The function starts by transferring the model to the specified compute device (either