

Learn to Perform Cubic Regression with Python: A Step-by-Step Guide

Authored by
Mohammed loot

April 27, 2026

RECOMMENDED CITATION

Mohammed loot (2026). *Learn to Perform Cubic Regression with Python: A Step-by-Step Guide*. PSYCHOLOGICAL STATISTICS. Retrieved from <https://statistics.arabpsychology.com/?p=3507>

Cubic regression represents a highly effective statistical methodology employed for modeling the relationship between a [predictor variable](#) and a [response variable](#), particularly when the underlying interaction exhibits a distinctive, complex [non-linear](#) structure. Distinct from the simplicity of linear or the single-curve nature of quadratic models, cubic regression possesses the unique capability to accurately capture trends that feature up to two turning points or inflection points. This modeling flexibility often results in characteristic 'S' or inverted 'S' shaped curves, making the technique indispensable for researchers and data analysts who must accurately interpret intricate, evolving data patterns across diverse domains, including engineering, biological sciences, and advanced economics.

This comprehensive, step-by-step guide is designed to systematically walk you through the entire process of performing [cubic regression](#) within the [Python](#) ecosystem. We will harness the collective power of crucial data science libraries: **pandas** for efficient data management, **NumPy** for robust numerical computation and model fitting, and **Matplotlib** for insightful data visualization. The tutorial covers every essential phase, starting with initial data preparation and exploratory visualization, progressing through the calculation and interpretation of the regression equation, and concluding with a rigorous evaluation of the model's performance using the widely accepted [R-squared](#) metric. Upon completion, you will possess the practical knowledge required to confidently apply this advanced [regression analysis](#) technique to your own datasets, thereby uncovering complex relationships that simpler statistical tools often overlook.

Understanding the Fundamentals of Cubic Regression

At its theoretical foundation, [cubic regression](#) is categorized as a specific type of polynomial regression, which extends the familiar concept of linear regression by allowing the relationship between variables to be mathematically represented by an n-th [degree](#) polynomial. Specifically, cubic regression mandates the use of a third-degree polynomial, meaning the standard regression equation adheres to the general mathematical form presented below:

$$y = \beta_0 + \beta_1x + \beta_2x^2 + \beta_3x^3 + \varepsilon$$

In this formula, y denotes the **response variable**, x represents the **predictor variable**, β_0 , β_1 , β_2 , β_3 are the regression coefficients that the model estimates from the data, and ε symbolizes the irreducible error term. The inclusion of the cubed term (x^3) grants the model significantly greater flexibility than lower-order models, enabling it to describe patterns where the rate of change is dynamic, changing direction multiple times across the range of the predictor.

The predominant advantage of employing [cubic regression](#) lies in its unparalleled ability to accurately model complex, curvilinear phenomena that linear or quadratic forms would struggle to represent. Consider, for example, studies in biological systems where population growth might follow a distinct cubic trajectory: initially accelerating rapidly, then slowing down (decelerating) as

resources become scarce or carrying capacity is approached, and finally accelerating once more due to evolutionary adaptations or the introduction of new resource vectors. Similarly, in financial markets, the relationship between advertising spending and total profit might exhibit such a nuanced trend. Identifying these intricate, multi-stage patterns is absolutely essential for generating accurate forecasts and formulating maximally informed strategic decisions, as the model explicitly captures the dual inflection points present in the data.

It is vital to clearly distinguish [cubic regression](#) within the broader spectrum of [regression analysis](#) techniques. While linear regression assumes a straight-line relationship (first [degree](#)) and quadratic regression models a single parabolic curve (second [degree](#)), cubic regression elevates the complexity by introducing the x^3 term, which allows for two distinct inflection points. This makes it ideally suited for datasets that demonstrate both increasing and decreasing trends across the observed range. However, this increased model complexity inherently raises the risk of **overfitting**--a scenario where the model fits the noise in the training data too closely, leading to poor generalization performance on new, unseen data. Therefore, the careful selection of the polynomial [degree](#) and robust evaluation procedures are critical to ensure the model's reliability and broad applicability.

Setting Up Your Python Environment

Before commencing the practical implementation of [cubic regression](#), it is a fundamental prerequisite to ensure that your [Python](#) environment is properly configured with the necessary libraries. This tutorial relies on three essential packages that form the backbone of nearly all data science and statistical computing workflows in [Python](#): **pandas**, **NumPy**, and **Matplotlib**. These libraries provide specialized, high-performance tools required for data handling, numerical operations, and visualization, respectively.

pandas: This library is indispensable for efficient data manipulation and structuring. It provides the powerful [DataFrame](#) object, which is the standard, high-performance structure used for managing tabular data.

NumPy: Serving as the foundational package for numerical computing, NumPy provides superior N-dimensional array objects and a comprehensive suite of mathematical functions optimized for array operations. It is the core tool we will use for the polynomial fitting itself.

Matplotlib: This versatile plotting library allows us to create static, animated, and interactive visualizations. It is crucial for generating the initial scatterplot and, subsequently, for visually overlaying the fitted regression curve onto the raw data points to assess model performance.

To begin, open your preferred [Python](#) development interface--whether it be a Jupyter Notebook, Spyder, or a simple Python script. If any of these libraries are not currently installed in your

environment, they can be easily added using pip, Python's package installer, by running the command `pip install pandas numpy matplotlib` in your terminal or command prompt. Once installed, you must import them into your active session using their conventional aliases to access their extensive functionalities efficiently:

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
```

These import statements establish the necessary environment, making the functionalities available through the concise aliases: `pd` for pandas, `np` for NumPy, and `plt` for Matplotlib's pyplot module. With this robust setup, we are ready to proceed to the crucial steps of data preparation, model implementation, and comprehensive visualization required for our cubic regression analysis.

Preparing and Visualizing Your Data

The initial and arguably most critical phase in any comprehensive [regression analysis](#) is the meticulous preparation of the data and the execution of exploratory visualization to gain crucial insights into its inherent structure and underlying patterns. For practical demonstration purposes in this guide, we will work with a controlled, hypothetical dataset composed of two variables: `x`, designated as our primary **predictor variable**, and `y`, serving as our **response variable**. This data will be structured and managed within a [pandas DataFrame](#), the standard tool for handling tabular data efficiently in Python.

The following code snippet demonstrates the process of creating this [DataFrame](#) directly from defined lists. It is considered a strong best practice to immediately print or display the resulting DataFrame after its creation. This step ensures rapid verification that the data has been loaded and structured without error, offering an instantaneous overview of the values and their respective ranges, which is paramount for early identification of potential data entry issues or unexpected distributions.

```
import pandas as pd
```

```
#create DataFrame
df = pd.DataFrame({'x': ,
'y': })
```

```
#view DataFrame
print(df)
```

```
x y
```

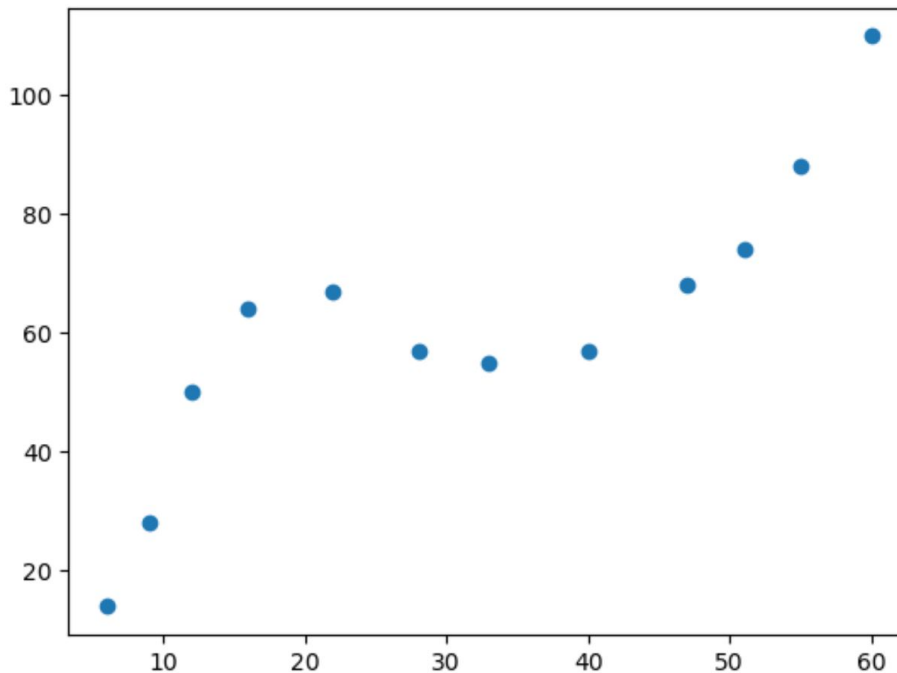
```
0 6 14
1 9 28
2 12 50
3 16 64
4 22 67
5 28 57
6 33 55
7 40 57
8 47 68
9 51 74
10 55 88
11 60 110
```

With the data successfully loaded, the next vital action is to visualize the relationship between the **predictor variable** x and the **response variable** y using a [scatterplot](#). This initial visualization is fundamentally crucial for identifying the true nature of the relationship--specifically, whether it follows a simple linear path or exhibits a more complex [non-linear](#) form. Visual assessment helps pre-determine the most appropriate modeling strategy, saving time and improving model accuracy.

```
import matplotlib.pyplot as plt
```

```
#create scatterplot
plt.scatter(df.x, df.y)
```

Upon reviewing the resulting [scatterplot](#), the relationship between the variables is clearly and distinctly [non-linear](#). As the value of x increases, the value of y initially rises, then appears to plateau and subsequently decline slightly, before finally rising sharply again. This specific pattern, defined by two discernible "turns" or inflection points within the data distribution, serves as powerful visual evidence that a [cubic relationship](#) is present. This compelling visual justification confirms that proceeding with a cubic regression model is the most appropriate analytical path, as it is uniquely capable of accurately capturing such complex, multi-directional fluctuations.



Fitting the Cubic Regression Model

With the [non-linear](#) nature of our data visually confirmed, the logical and necessary next stage is the fitting of the [cubic regression](#) model. For implementation within the [Python](#) environment, we will utilize the robust numerical tools provided by the **NumPy** library. Specifically, two key functions are employed: `numpy.polyfit()`, which is used to calculate the optimal polynomial coefficients that minimize the squared error, and `numpy.poly1d()`, which then constructs a convenient, functional polynomial object from those derived coefficients.

The core function, `numpy.polyfit()`, requires three primary inputs: the x-coordinates (our **predictor variable**), the y-coordinates (our **response variable**), and the desired [degree](#) of the polynomial. Since we are performing cubic regression, it is essential to specify a [degree](#) of 3. This function returns an array of the estimated coefficients, ordered descendingly from the highest [degree](#) term down to the intercept (constant) term. These coefficients are then passed to `numpy.poly1d()`, which creates a callable object. This object behaves precisely like the mathematical function $y = f(x)$, allowing us to easily compute predicted y-values for any given x-values, which is crucial for plotting the fitted line.

The following code snippet demonstrates the practical implementation of fitting the [cubic regression](#) model and subsequently visualizing the resulting fitted regression curve by overlaying it onto our original [scatterplot](#). To ensure the fitted line appears smooth and continuous across the data range, we generate a dense sequence of evenly spaced x-values using `numpy.linspace()`. This visualization step is critically important as it provides an immediate, intuitive assessment of

how accurately the estimated cubic model aligns with the observed data points.

import numpy as np

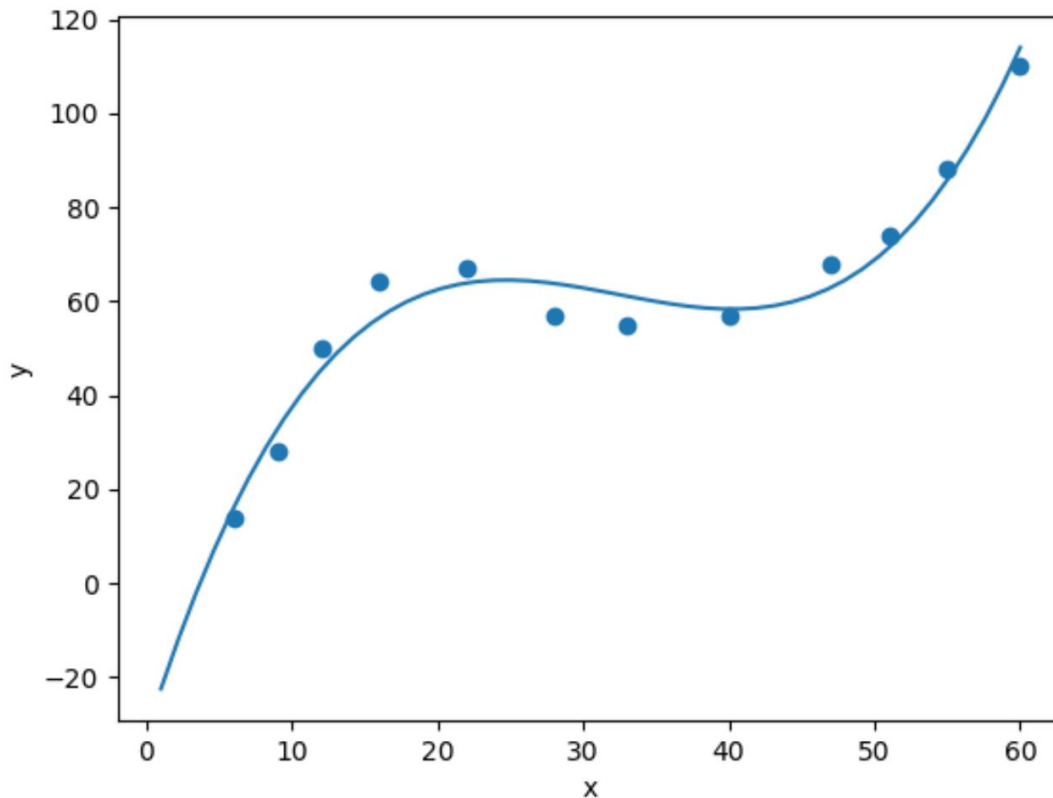
```
#fit cubic regression model
model = np.poly1d(np.polyfit(df.x, df.y, 3))

#add fitted cubic regression line to scatterplot
polyline = np.linspace(1, 60, 50)
plt.scatter(df.x, df.y)
plt.plot(polyline, model(polyline))

#add axis labels
plt.xlabel('x')
plt.ylabel('y')

#display plot
plt.show()
```

The resulting plot below provides strong visual confirmation that the [cubic regression](#) line successfully captures the complex [non-linear](#) trend observed in our dataset. The fitted curve accurately reflects the initial sharp increase, the subsequent period of decline and stabilization, and the final resurgence of the **response variable** \bar{y} as the **predictor variable** \bar{x} progresses. This close visual alignment significantly increases our confidence in the model's ability to represent the true data generating process.



To precisely quantify the mathematical relationship established by our fitted model, we can simply print the `model` object, which was created using `numpy.poly1d()`. This action directly outputs the estimated [polynomial regression](#) equation, detailing the numerical values of the coefficients (β_3 , β_2 , β_1 , β_0) that define the specific cubic curve best suited for our observed data.

print(model)

```
3 2
0.003302 x - 0.3214 x + 9.832 x - 32.01
```

Based on this output, we can formally articulate the final fitted [cubic regression](#) equation as:

$$\hat{y} = 0.003302x^3 - 0.3214x^2 + 9.832x - 32.01$$

This powerful equation is now immediately usable for making predictions. We can substitute any value of \hat{x} , particularly within the observed range of our data, to estimate the corresponding expected value of \hat{y} . For instance, if we wish to predict the response when the predictor \hat{x} equals 30, the calculation confirms the predictive utility of the model:

$$\hat{y} = 0.003302(30)^3 - 0.3214(30)^2 + 9.832(30) - 32.01 = 64.844$$

This demonstrates the model's effectiveness in quantifying and forecasting based on the established [non-linear](#) relationship.

Evaluating Model Performance with R-squared

After successfully fitting a [cubic regression](#) model, the next mandatory step is to rigorously evaluate its performance and quantify how effectively it accounts for the inherent [variance](#) in the **response variable**. The most common and widely accepted measure for this purpose is the [R-squared](#) (R^2) statistic, officially known as the coefficient of determination. [R-squared](#) provides a value, typically between 0 and 1, that measures the proportion of the total [variance](#) in the dependent variable that is predictable from the independent variable(s) included in the regression model.

A higher [R-squared](#) value, approaching 1 (or 100%), indicates a superior fit, meaning that a very large proportion of the variability in the **response variable** \hat{y} is successfully explained by the fluctuations in the **predictor variable** x , suggesting strong predictive power. Conversely, a low [R-squared](#) value suggests the model is inadequate, explaining only a small fraction of the total variability, and implies that significant influential factors are missing or that a different [regression analysis](#) technique might be required. For polynomial models, the calculation of [R-squared](#) requires comparing the explained [variance](#) (Sum of Squares Regression) to the total variance (Total Sum of Squares), effectively measuring the improvement gained over simply using the mean to predict the response.

While **NumPy**'s [polyfit\(\)](#) and [poly1d\(\)](#) functions are exceptionally efficient for fitting polynomial models, they unfortunately do not provide the [R-squared](#) statistic directly as part of their standard output. Therefore, to obtain this critical evaluation metric, we must define a concise, custom function. This function encapsulates the entire process: it performs the polynomial fit, computes the predicted values (\hat{y}), and then calculates the necessary sum of squares components to accurately derive the final [R-squared](#) value, thereby providing a quantitative measure of the explained variance for any specified [degree](#).

#define function to calculate r-squared

```
def polyfit(x, y, degree):  
    results = {}  
    coeffs = np.polyfit(x, y, degree)  
    p = np.poly1d(coeffs)  
    #calculate r-squared  
    yhat = p(x)  
    ybar = np.sum(y)/len(y)  
    ssreg = np.sum((yhat-ybar)**2)
```

```
sstot = np.sum((y - ybar)**2)
results = ssreg / sstot

return results

#find r-squared of polynomial model with degree = 3
polyfit(df.x, df.y, 3)
```

Execution of this function, specifying a [degree](#) of 3 for our dataset, yields the following output, which represents the [R-squared](#) value for our fitted [cubic regression](#) model:

```
{'r_squared': 0.9632469890057967}
```

In this specific instance, the calculated [R-squared](#) value is approximately **0.9632**. This exceptionally high score carries the significant interpretation that 96.32% of the total variation observed in the **response variable** (\bar{y}) is effectively explained and accounted for by the **predictor variable** (x) through the fitted [cubic regression](#) model. Such a robust result confirms that our model provides an excellent fit to the data, demonstrating its strong capability to accurately capture the intricate [non-linear](#) relationship that was initially observed in the exploratory scatterplot visualization.

Conclusion and Further Exploration

This tutorial successfully navigated the entire process of performing [cubic regression](#) in [Python](#), guiding the reader through every essential analytical phase. We began with a theoretical understanding of the technique and its critical role in modeling complex [non-linear](#) patterns. We then proceeded through the practical setup of the [Python](#) environment, the necessary steps of data preparation and visualization, the technical fitting of the cubic model using **NumPy**, the mathematical interpretation of the derived regression equation, and finally, the rigorous evaluation of the model's goodness of fit via the [R-squared](#) metric.

The proficiency to accurately model [non-linear](#) relationships, especially those characterized by multiple turning points or inflection curves, is an absolutely essential skill set in advanced data analysis across nearly all scientific and business applications. [Cubic regression](#) stands out as a particularly effective and versatile tool for addressing scenarios where simpler linear or quadratic models are fundamentally incapable of capturing the nuanced dynamics inherent in the data. By following the detailed, systematic methodology presented in this guide, you are now equipped to confidently apply this sophisticated technique to diverse datasets, enabling you to extract deeper, more meaningful insights into complex inter-variable relationships and formulate substantially more accurate and reliable predictions.

While the [R-squared](#) value offers a highly valuable and intuitive measure of model fit, it is important to emphasize that it is only one component of a comprehensive model evaluation strategy. In professional analytical contexts, it is always advisable to consider additional diagnostic metrics, such as the adjusted R-squared (which penalizes the inclusion of excessive predictors), the Root Mean Squared Error (RMSE), and visual qualitative assessments derived from residual plots. Furthermore, always strive to select the optimal model complexity--specifically, the [degree](#) of the polynomial--that skillfully balances maximizing explanatory power against mitigating the inherent risk of overfitting the training data, ensuring the model's generalizability and reliability on future observations.

Additional Resources

To further enhance your understanding of [regression analysis](#) and related data science topics in [Python](#), we highly recommend exploring the following tutorials and official documentation sources:

[How to Perform Polynomial Regression in Python](#)

[NumPy polyfit\(\) Official Documentation](#)

[NumPy poly1d\(\) Official Documentation](#)

[Pandas Official Documentation](#)

[Matplotlib Official Documentation](#)