

Multiple Regression(with Polynomial Regression)

Muhammad Moaz Amin
dept. Electronic Engineering
Hochschule Hamm Lippstadt
Hamm, Germany
muhammad-moaz.amin@stud.hshl.de

Abstract—Estimating the relationship among different variable which have reason and result relation holds great importance, to estimate that Regression analysis is a statistical technique that is adapted. Main focus of this uni-variate regression is to analyse the relationship between a non linear and a linear variable and to formulate a linear equation between the two. a regression model which contain one linear and multiple non-linear independent variables is most often called multi linear regression. [3] This paper is concentrated on the polynomial regression model, which is useful when there is reason to believe that relationship between two variables is curvilinear. The polynomial regression model has been applied using the characterisation of the connection between strains and drilling depth. Parameters of the model were estimated employing a least square method. After fitting, the model was evaluated using a number of the common indicators wont to evaluate accuracy of regression model.

I. INTRODUCTION

Regression analysis involves identifying the link between a variable quantity and one or more independent variables. It's one amongst the foremost important statistical tools which is extensively utilized in the majority sciences. It's specially used in business and economics to check the link between two or more variables that are related causally. A model of the relationship is hypothesized, and estimates of the parameter values are accustomed develop an estimated equation.

Common questions which are generally asked in this research of Multiple regression analysis generally revolve around "are there any relations between dependent and independent variable?", and "if there are some relations that exist, what is total power of the relation?," "is there any possibility to predict about orientation regarding the dependent variable?", and "if certain conditions are controlled, what influences does a special variable or a group of variables have over another variable or variables?". [1]

Uni-variate analysis is the regression which uses a single independent variable while multivariate regression analysis uses more than one independent variable [8] [5]. The relation between dependent variable and an independent variable is analysed through uni-variate regression analysis, and the equation which represents the linear relationship between the both is formulated.

In Multivariate regression analysis, an attempt is made to account for the variation between the dependent variable and

the independent variable synchronically [2]. Formulation of Multivariate model is represented in Figure 1.

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_n x_n + \epsilon$$

y - dependent variable
 X_i - independent variable
 β_i - parameter
 ϵ - error

Fig. 1. Multivariate Model

The assumptions of multivariate regression analysis are normal distribution, linearity, freedom from extreme values and having no multiple ties between independent variable [5].

II. MULTIPLE REGRESSION WITH PYTHON

Linear Regression generally is a supervised method for machine learning rooted in statistics. From this method numeric values are forecasted using a combination of predictors which can be both numeric or binary variables. The condition for getting the variable depends on a certain relation at hand (a linear, measurable by a correlation) with the target variable. [4]

However, in reality there are a number of factors which alter the results of the working predictive model. There are usually more than variables that work together to achieve better and reliable results from a prediction. This causes more complexity in our model and hence representing it on a two- dimensional plot is not easy. All the predictors will constitute their own unique dimension and we would have to assume that our predictors apart from being related to the response are also related among themselves and this characteristic of data is called multicollinearity. [4]

A. Multiple Regression Formulation

The basic Multiple Regression is described in Figure 1 and in depth detail is shown in Figure 2 where dependent(response) variable Y on a set of k independent(predictor) variables X_1, X_2, \dots, X_k can be expressed as where

- y_i = value of dependent variable, Y is for i th case.
- x_{ij} = value of j th independent variable, X_j for i th case.
- β_0 is the Y -intercept of the regression surface.

$$\begin{cases} y_1 = \beta_0 + \beta_1 x_{11} + \dots + \beta_k x_{1k} + e_1 \\ y_2 = \beta_0 + \beta_1 x_{21} + \dots + \beta_k x_{2k} + e_2 \\ \vdots \\ y_n = \beta_0 + \beta_1 x_{n1} + \dots + \beta_k x_{nk} + e_n \end{cases}$$

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_k x_{ik} + e_i, \text{ for } i = 1, 2, \dots, n$$

Fig. 2. In depth Multiple Regression Model showing relation of independent(k) and dependent variable(Y)

- each $\beta_j, j=1,2,\dots,k$, is the slope of the regression surface w.r.t. variable X_j and e_i is the random error component for the i th case.

In the first equation we have n observations and k predictors ($n > k+1$) The assumptions of the multiple regression model are similar to those for the simple linear regression model. Model assumptions [4]:

B. Observations

- errors e_i are normally distributed with their mean as zero and their standard deviation σ are independent of the error terms associated with all other observations. Errors are not related to each other.
- variables X_j in the context of regression analysis are considered as fixed quantities, whereas they are random variables in the context of correlation analysis. But in both the cases X_j are totally independent of the error term. If X_j are assumed as fixed quantities, then we are assuming that we have realizations of k variables X_j and the only randomness in Y is coming from the error term.

In matrix form, we can rewrite the regression model as described in Figure. 3.

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$$

Fig. 3. Matrix notation of Model

- response vector \mathbf{Y} and error vector \mathbf{e} are the column vectors of length n .
- vector of parameters $\boldsymbol{\beta}$ is the column vector of length $k+1$ and its design matrix (where all elements in the first column are equal to 1, and second column's values are filled by the observed values of X_1 , etc).

Here, values of $\boldsymbol{\beta}$ and \mathbf{e} are unknown and assumed.

III. POLYNOMIAL REGRESSION

A. History

The first design of an experiment for polynomial Regression appeared in a paper in 1815 by Gergonne. [7] In the twentieth century, polynomial regression played an important role in the development of regression analysis with greater emphasis on issues of design and interface. [6]. Recently, polynomial models have been complemented by other methods, with

no-polynomial models being more advantageous for certain classes of problems.

B. Definition

Polynomial Regression is a technique of Multiple Regression which provides an automatic means of creating both interactions and non-linear power transformations of the original variables systematically. The number of bends that fit the curve depend on the degree of power. Higher the power higher will be the bends.

For Example in simple linear regression form:

$$y = \beta_0 + \beta_1 x$$

With a second degree transformation of the simple form, **quadratic** form is achieved:

$$y = \beta_0 + \beta_1 x + \beta_2 x^2$$

Likewise, with a third degree transformation the equation turns into **cubic**:

$$y = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3$$

In case of regression being a multiple one, the expansion will create more terms as the power is increasing resulting in more features being derived. A regression which is a result of two predictors (x_1 and x_2 which is expanded using quadratic transformation will result in the following equation:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2 + \beta_4 x_1^2 + \beta_5 x_2^2$$

With such an expansion two important aspects are observed:

- Polynomial expansion increases number of predictors rapidly.
- Power of predictors is directly proportional to degree of polynomials, which results in deteriorating numeric stability, thus it requires standardized numeric values which are too large or suitable numeric formats.

Matrix Form of Polynomial Regression is described in Figure 5. [3]

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_1 & x_1^2 & \dots & x_1^m \\ 1 & x_2 & x_2^2 & \dots & x_2^m \\ 1 & x_3 & x_3^2 & \dots & x_3^m \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & x_n^2 & \dots & x_n^m \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_m \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \epsilon_3 \\ \vdots \\ \epsilon_n \end{bmatrix},$$

Fig. 4. Matrix Model of Polynomial Regression

C. Why do we need Polynomial Regression?

In the case of Simple Linear Regression the best fit formed line is a straight line where as actual values form a curve which results in formed line not cutting the mean of the points.

In such cases, Polynomial Regression is beneficial, as it predicts the best fit line that follows the pattern of the plotted points resulting in forming a curve. This is described in the Figure 4.

One of the few differences between Polynomial Regression and Linear Regression is that Polynomial Regression does not require the relationship between the independent and dependent variables to be linear in the data set. Moreover, Polynomial Regression Model is used when the Linear Regression fails to describe the best result and the points are not captured by it.

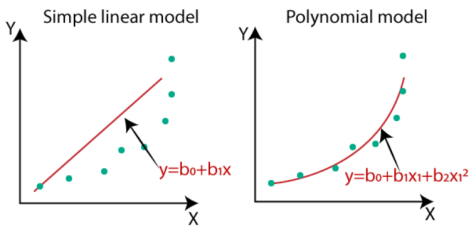


Fig. 5. Linear and Polynomial Regression plotted

IV. WHY SHOULD WE CHOOSE POLYNOMIAL REGRESSION?

Linear Regression requires the relation between the dependent and the independent variable to be linear. But in the case of distribution data being more complex? Is it possible to fit non-linear data into linear models? How can a curve be produced that captures the data as in Figure 6? To understand

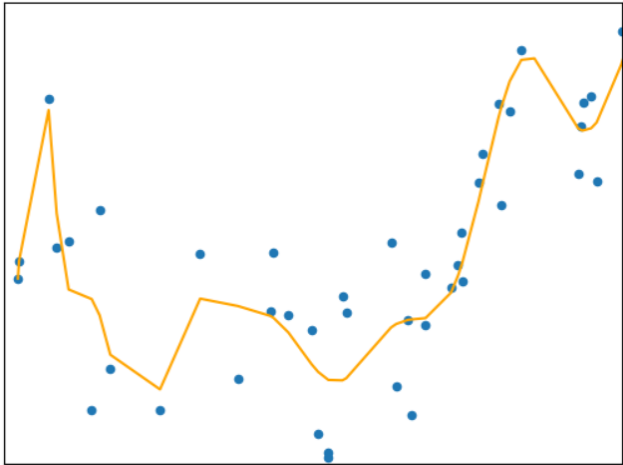


Fig. 6. Complex data Set

something’s importance, it’s comparison with its competitor is very important. Just like that to understand the need of Polynomial Regression we need to first evaluate set of data

with linear Regression and then produce it using Polynomial Regression. at First we need to generate some random data:

```
import numpy as np
import matplotlib.pyplot as plt
np.random.seed(0)
x = 2- 3 * np.random.normal(0,1,20)
y = x - 2 * (x ** 2) + 0.5 * (x ** 3) + np.random.normal(-3,
3, 20)
plt.scatter(x,y, s=10)
plt.show()
```

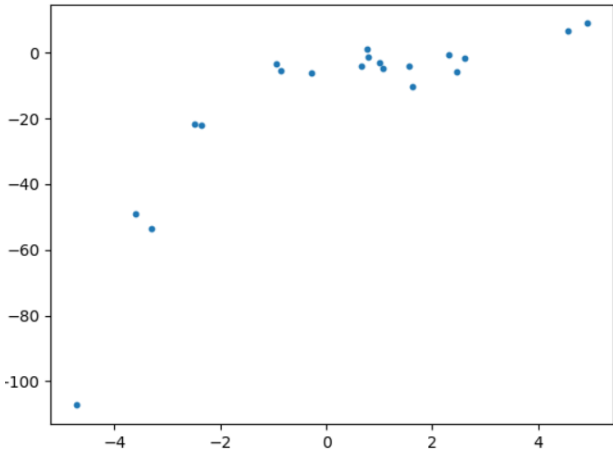


Fig. 7.

TABLE I
TABLE TYPE STYLES

Table Head	Table Column Head		
	Table column subhead	Subhead	Subhead
copy	More table copy ^a		

^aSample of a Table footnote.



Fig. 8. Example of a figure caption.

Figure Labels: Use 8 point Times New Roman for Figure labels. Use words rather than symbols or abbreviations when writing Figure axis labels to avoid confusing the reader. As an example, write the quantity “Magnetization”, or “Magnetization, M”, not just “M”. If including units in the label, present them within parentheses. Do not label axes only with units. In the example, write “Magnetization (A/m)” or “Magnetization {A[m(1)]}”, not just “A/m”. Do not label axes with a ratio of quantities and units. For example, write “Temperature (K)”, not “Temperature/K”.

REFERENCES

- [1] Pandora - uygulamalı Çok deęişkenli İstatistiksel yöntemlere giriş 1 - reha alpar - kitap - isbn 9789755914312.
- [2] Ünver. Ö. Gamgam. Uygulamalı İstatistik i - İstanbul.
- [3] Güliden Kaya Uyanık and Neşe Güler. A study on multiple linear regression analysis. *Procedia - Social and Behavioral Sciences*, 106:234–240, 12 2013.
- [4] Luca Massaron. Regression analysis with python: Nook book, Feb 2016.
- [5] Pegem.Net. Sosyal bilimler için veri analizi el kitabı İstatistik, araştırma deseni spss uygulamaları ve yorum - Şener büyüköztürknbps:: Kpss, Öabt, ales, dgs, yks, lgs, yds, gys kitapları: Pegem.net İnternetteki kitapçınız.
- [6] Kirstine Smith. On the Standard Deviations of Adjusted and Interpolated Values of an Observed Polynomial Function and its Constants and the Guidance they give Towards a Proper Choice of the Distribution of Observations, November 1918.
- [7] Stephen M Stigler. Gergonne’s 1815 paper on the design and analysis of polynomial regression experiments. *Historia Mathematica*, 1(4):431–439, 1974.
- [8] Barbara G. Tabachnick and Linda S. Fidell. *Using Multivariate Statistics (5th Edition)*. Allyn Bacon, Inc., USA, 2006.

IEEE conference templates contain guidance text for composing and formatting conference papers. Please ensure that all template text is removed from your conference paper prior to submission to the conference. Failure to remove the template text from your paper may result in your paper not being published.