

REGRESSION ANALYSIS AND APPROXIMATION
BY MEANS OF CHEBYSHEV POLYNOMIAL
REGRESIVNA ANALIZA I APROKSIMACIJA POMOĆU
ČEBIŠEVLJEVIH POLINOMA

Nikola Tomašević, Marko Tomašević, Tatjana Stanivuk*

Faculty of Economics, University of Belgrade, Belgrade, Serbia*; Faculty of Maritime Studies, University of Split, Split, Croatia
Ekonomski fakultet, Sveučilište u Beogradu, Beograd, Srbija*; Pomorski fakultet, Sveučilište u Splitu, Split, Hrvatska

Abstract

The paper deals with the regression model, describes the procedure of getting regression coefficients and gives the analysis of the model variance. Chebyshev polynomials and their properties are also stated. The transformation of the regression model, from segment [a, b] to segment [-1, 1] is performed, as well as the approximation of the obtained regression polynomial, using prespecified accuracy polynomials of lower degree.

Sažetak

Rad se bavi regresivnim modelom, opisuje postupak dobivanja regresivnih koeficijenata, te daje analitički model varijance. Promatraju se i Čebiševljevi polinomi i njihova svojstva. Izvodi se transformacija regresivnog modela od segmenta [a, b] do segmenta [-1, 1], kao i aproksimacija izvedene regresije polinoma, koristeći unaprijed određene polinome nižeg stupnja.

Introduction

Trend, as a notion, implies the expected development (or a development tendency, as it is frequently called) of a phenomenon in a future period, based on the observation of that phenomenon over a previous and long enough period. More precisely, on the basis of phenomenon development data in the previous period, i.e. a series of data in previous consecutive periods, a mathematical equation, best corresponding to the data given, is selected. The equation is then "extended" to future periods, assuming that the phenomenon will be developing in the same way, and its future values are calculated. The values usually refer to a sufficiently close future period.

Most frequently, the selected functions for expressing regression are the following:

- The linear regression in the form of:

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i, \quad i = 1, 2, \dots, n$$

- Regression second degree polynomial (parabolic regression) in the form of:

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i + \hat{\beta}_2 x_i^2, \quad i = 1, 2, \dots, n$$

- Regression k degree polynomial ($\beta_0, \beta_1, \dots, \beta_k$) in the form of:

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i + \hat{\beta}_2 x_i^2 + \dots + \hat{\beta}_k x_i^k, \quad i = 1, 2, \dots, n$$

The calculation of function parameter values, which are used to express regression, is performed on the basis of the data already known from the previous period. The technique of determination of the coefficients itself is performed using the method of least squares. The parameters of the required function are selected to make the sum of square difference of known real values and the values of the required function minimum. We determine the value of the required coefficients by minimum, from these equations.

Regression analysis

Mutual connections among changeable phenomena can be divided into two groups: deterministic and stochastic. A deterministic connection, which is also called functional, refers to the case in which only one value of the independent

variable X corresponds to only one value of the dependent variable Y . This type of connection can be shown in a general expression: $Y = g(X)$, where $g(X)$ is a function of X . A stochastic connection refers to the case in which whole series of possible values of the dependent variable correspond to only one value of the independent variable. The essence of the stochastic connections to be observed in this paper is the functional connection which exists between individual values of the independent variable X and average values of the dependant variable Y (i.e. expected values).

Thus:

$$E[Y] = f(X)$$

The main problem of regression is to estimate the functional $(x_1, Y_1), \dots, (x_n, Y_n)$ dependence on the basis of sample, as well as to the give statistic evaluation of accuracy of such estimations.

The regression $Y = Y(x)$ model consists of two additive parts: deterministic and stochastic. Hence, the regression model can be observed as the following sum:

$$Y = g(X) + e_i, i = 1, 2, \dots, n$$

The specification of the regression model as a stochastic model does not only imply its mathematical expression. It also implies assumptions which ensure the optimal estimation of parameters. The following assumptions are most frequently introduced:

1. Linearity – There is a linear connection between individual values of the independent variable, X, x_i and the corresponding Y value, and also the average value $E[Y_i]$.
2. $E(e_i) = 0$
3. Homogeneity of variance – Stochastic members have the equal deviation, or more precisely, equal variances. Hence:

$$V[e_1] = V[e_2] = \dots = V[e_n] = \sigma^2.$$

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \dots \\ y_n \end{bmatrix} \quad X = \begin{bmatrix} 1 & x_1 & x_1^2 & \dots & x_1^k \\ 1 & x_2 & x_2^2 & \dots & x_2^k \\ \dots & \dots & \dots & \dots & \dots \\ 1 & x_n & x_n^2 & \dots & x_n^k \end{bmatrix} \quad \beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \dots \\ \beta_k \end{bmatrix} \quad e = \begin{bmatrix} e_1 \\ e_2 \\ \dots \\ e_n \end{bmatrix}$$

In this notation, Y stands for the vector of the empirical values of the dependent variable, X stands for the matrix of the value of the regression variable

4. There is no autocorrelation between $e_i, i = 1, 2, \dots, n$, hence, there is no linear connection between any of two stochastic members e_i and e_j .
5. X is not a random variable (therefrom, the independent variable in the model is denoted with the small letter). This assumption indicates that the values of the independent variable are fixed. In other words, they must be chosen in advance by the researcher before sample taking.
6. e_i has a normal distribution so it can be written as: , i.e. the stochastic member has a normal distribution with the arithmetic mean $e_i \sim \rightarrow (0, \sigma^2)$ equivalent to zero and σ^2 variance.

Our aim is to find the best possible coefficient estimations of the regression model on the basis of the sample. The method of finding coefficient estimations of the regression model on the basis of the sample, $(x_1, Y_1), \dots, (x_n, Y_n)$ which will be applied here, is the method of least squares.

The basic regression model is the k degree polynomial in the form of:

$$Y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \dots + \beta_k x_i^k + e_i, i = 1, 2, \dots, n$$

where:

- Y_i – the i dependant variable (project value)
- x_i – i value of the independent variable
- $\beta_0, \beta_1, \dots, \beta_k$ – regression parameters (regression coefficient)
- e_i – stochastic member or random error
- n – sample size

The system of equations in its matrix notation is

$$Y = X\beta + e$$

where the following has been set:

X, β stands for the vector of the unknown parameter values, and e stands for the vector of the unknown values of the random variable e_i .

$$E[e_i] = 0, \text{ per each } i$$

On the basis of the sample, the best possible estimations $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_k$ with the corresponding regression coefficients $\beta_0, \beta_1, \dots, \beta_k$ have to be found in order to determine the regression curve of the sample:

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i + \hat{\beta}_2 x_i^2 + \dots + \hat{\beta}_k x_i^k, i = 1, 2, \dots, n$$

In this sample the sign \hat{Y}_i denotes Y value which is positioned on the best adjusted regression curve and therefore is called the adjusted value of Y .

Thus, the minimum of the expression

$$\sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n \left[y_i - \left(\hat{\beta}_0 + \hat{\beta}_1 x_i + \hat{\beta}_2 x_i^2 + \dots + \hat{\beta}_k x_i^k \right) \right]^2 \text{ has to be found.}$$

In the above expression $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_k$ are unknown. The procedure for minimizing is performed by finding the partial derivatives of the expression above

As a rule, the regression curve of the basic set differs from that of the sample because the estimated values $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_k$ differ from the real parameter values. The vertical deviation $\beta_0, \beta_1, \dots, \beta_k$ (the difference) between the real value y_i and the adjusted value \hat{y}_i is called the residual value and it is denoted by e_i :

$$e_i = y_i - \hat{y}_i = y_i - \left(\hat{\beta}_0 + \hat{\beta}_1 x_i + \hat{\beta}_2 x_i^2 + \dots + \hat{\beta}_k x_i^k \right)$$

The main idea of the method of least squares is to choose the regression curve with the least sum of squares residue out of all regression curves possible.

by ??? and their equalization to zero. In this way, the following system of $k + 1$ equations is obtained. These equations are called normal equations:

$$\begin{aligned} n\hat{\beta}_0 + \hat{\beta}_1 \sum_{i=1}^n x_i + \dots + \hat{\beta}_k \sum_{i=1}^n x_i^k &= \sum_{i=1}^n y_i \\ \hat{\beta}_0 \sum_{i=1}^n x_i + \hat{\beta}_1 \sum_{i=1}^n x_i^2 + \dots + \hat{\beta}_k \sum_{i=1}^n x_i^{k+1} &= \sum_{i=1}^n x_i y_i \\ \dots & \dots \\ \hat{\beta}_0 \sum_{i=1}^n x_i^k + \hat{\beta}_1 \sum_{i=1}^n x_i^{k+1} + \dots + \hat{\beta}_k \sum_{i=1}^n x_i^{2k} &= \sum_{i=1}^n x_i^k y_i \end{aligned}$$

Their matrix notation is the following:

$$(X^T X) \cdot \hat{\beta} = X^T Y, \quad \text{i.e.} \quad \hat{\beta} = (X^T X)^{-1} \cdot (X^T Y),$$

where

$$X^T X = \begin{bmatrix} n & \sum x_i & \dots & \sum x_i^k \\ \sum x_i & \sum x_i^2 & \dots & \sum x_i^{k+1} \\ \dots & \dots & \dots & \dots \\ \sum x_i^k & \sum x_i^{k+1} & \dots & \sum x_i^{2k} \end{bmatrix}; \quad X^T Y = \begin{bmatrix} \sum y_i \\ \sum x_i y_i \\ \dots \\ \sum x_i^k y_i \end{bmatrix}.$$

If the assumptions of the regression model are fulfilled, the estimations obtained by the method of least squares will be the best, i.e. the most efficient

and impartial linear estimations. Thus, the following is valid:

$$E[\hat{\beta}_0] = \beta_0, \quad E[\hat{\beta}_1] = \beta_1, \dots, E[\hat{\beta}_k] = \beta_k,$$

i.e. the estimations $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_k$ are equal, on average, to the unknown parameters with respect to $\beta_0, \beta_1, \dots, \beta_k$. The standard deviation is denoted by $\sigma_{\hat{\beta}_i}$ and it is called the standard error of the estimation of $\hat{\beta}_i$. It measures the estimation of $\hat{\beta}_i$ deviation from parameters β_i ($i=1, \dots, k$).

In order to consider the components of variability of the dependant variable, the sample regression drawn in the diagram (as in Figure 1) has to be observed. The diagram shows an arbitrary empirical value y_i taken from the sample which corresponds to the value of the dependant variable x_i .

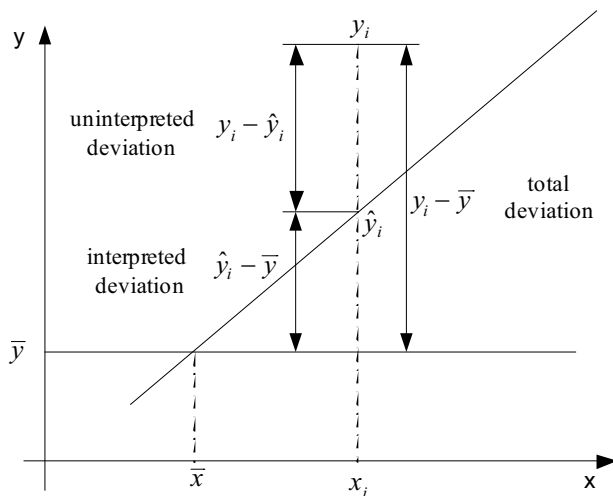


Figure 1

The total deviation of the dependant variable can be treated as the sum of interpreted and uninterpreted deviation:

$$(y_i - \bar{y}) = (\hat{y}_i - \bar{y}) + (y_i - \hat{y}_i).$$

$y_i - \bar{y}$ - total deviation

$\hat{y}_i - \bar{y}$ - interpreted deviation

$y_i - \hat{y}_i$ - uninterpreted deviation

This equality will be valid if we square and sum both sides by all values of the sample, so it can be said that the total variability equals to the sum of interpreted and uninterpreted variability:

$$e_i = y_i - \hat{y}_i = y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i + \hat{\beta}_2 x_i^2 + \dots + \hat{\beta}_k x_i^k)$$

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^n (\hat{y}_i - y_i)^2$$

i.e. $ST = SP + SR$

ST - total sum of squares

SP - interpreted sum of squares

SR - uninterpreted sum of squares

It is known that the sample variance is obtained when the sum of squares is divided by the number of degrees of freedom ($n - k - 1$). The estimated variance of regression model is obtained by the following formula:

$$s^2 = \hat{\sigma}^2 = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n - (k + 1)},$$

or in matrix notation:

$$\hat{\sigma}^2 = \frac{Y^T Y - \hat{\beta}^T (X^T Y)}{n - (k + 1)}$$

The standard regression error is the square root of the residual variance and denotes the standard deviation of random error:

$$s = \hat{\sigma} = \sqrt{\frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n - (k + 1)}} = \sqrt{\frac{Y^T Y - \hat{\beta}^T (X^T Y)}{n - (k + 1)}}.$$

The standard error of the regression model parameters is given by the formula:

$$se(\hat{\beta}_j) = \hat{\sigma} \sqrt{s_{jj}}, j = 1, 2, \dots, k$$

where s_{jj} stands for diagonal matrix $(X^T X)^{-1}$ elements.

The following text will present the procedure for obtaining individual elements in the equation decomposition of sum of squares of deviation of empirical values of the dependant variable from the average value of the same variable.

Starting from the following formula:

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (y_i^2 - 2\bar{y}y_i + \bar{y}^2) = \sum_{i=1}^n y_i^2 - n\bar{y}^2,$$

that is:

$$ST = Y^T Y - n\bar{Y}^2;$$

$$\sum_{i=1}^n (\hat{y}_i - \bar{y})^2 = \hat{\beta}_0 \sum_{i=1}^n y_i + \hat{\beta}_1 \sum_{i=1}^n x_i y_i + \dots + \hat{\beta}_k \sum_{i=1}^n x_i^k y_i - n\bar{y}^2,$$

i.e.

$$SP = \hat{\beta}^T (X^T Y) - n\bar{Y}^2,$$

hence:

$$SR = Y^T Y - \hat{\beta}^T (X^T Y)$$

The respective number of degrees of freedom with the total sum of squares is , with the interpreted part k and the residual sum of squares $n - (k + 1)$.

The method of variance analysis examines whether the regression parameter of, for example, k - degree reduces significantly the residual sum of squares. In other words, it examines if the k -degree polynomial would be optimal.

Approximation using Chebyshev polynomials

In the analysis of the variance, when finding the optimal degree of the polynomial and testing its coefficient, it is very important for the polynomial to be of the least degree possible, because the mathematical calculation and search for theoretical data using the optimal polynomial may result in errors. These errors might accumulate and turn into system errors, which may lead to wrong conclusions. This has been the reason for choosing Chebyshev polynomials to reduce the polynomial degree and accelerate the process of search for theoretical data, keeping at the same time the coefficients given within the estimated intervals. In finding solutions to this problem, some areas of special functions have been used and they will be stated in the following pages.

Namely, it is not difficult to prove that, using the linear transformation

$$x = \frac{b-a}{2}t + \frac{b+a}{2},$$

polynomials

$$\tilde{T}_m(x) = \left(\frac{b-a}{2}\right)^m T_m\left(\frac{x - \frac{b+a}{2}}{\frac{b-a}{2}}\right)$$

keep at the segment $[a,b]$ all properties of Chebyshev polynomials $T_m(t)$ defined at the segment $[-1,1]$. It is also possible to prove that:

$$\max_{a \leq x \leq b} (\tilde{T}_m(x) - 0) = \tilde{E}_m = \left(\frac{b-a}{2}\right)^m \cdot \frac{1}{2^{m-1}} = 2 \left(\frac{b-a}{4}\right)^m$$

Since

$$T_n(x) = \cos(n \arccos x),$$

and, if starting from the identity

$$\cos(n+1)a + \cos(n-1)a = 2 \cos na - \cos a$$

we put $a = \arccos x$, we get :

$$T_{n+1}(x) = 2xT_n(x) - T_{n-1}(x) \quad (-1 \leq x \leq 1)$$

This relation can be applied to Chebyshev polynomial formation. Indeed, starting from

$$T_0(x) = 1 \text{ and } T_1(x) = x, \text{ and using (1) we get:}$$

$$T_2(x) = 2x \cdot T_1 - T_0 = 2x^2 - 1,$$

$$T_3(x) = 2x \cdot T_2 - T_1 = 4x^3 - 3x,$$

$$T_4(x) = 2x \cdot T_3 - T_2 = 8x^4 - 8x^2 + 1$$

This procedure can be further extended. Using the above formulas, we shall approximate the polynomial:

$$f(x) = a_0 + a_1x + a_2x^2 + a_3x^3 + a_4x^4$$

by the third degree polynomial, so that:

$$\max_{-1 \leq x \leq 1} |f(x) - P_3(x)| < 0,05.$$

Using the expressions above we have:

$$1 = T_0, \quad x = T_1, \quad x^2 = \frac{1}{2}(T_0 + T_2)$$

$$x^3 = \frac{1}{4}(3T_1 + T_3), \quad x^4 = \frac{1}{8}(3T_0 + 4T_2 + T_4),$$

hence:

$$f(x) = a_0T_0 + a_1T_1 + \frac{a_2}{2}(T_0 + T_2) + \frac{a_3}{4}(3T_1 + T_3) + \frac{a_4}{8}(3T_0 + 4T_2 + T_4) =$$

$$= \left(a_0 + \frac{a_2}{2} + \frac{3a_4}{8}\right)T_0 + \left(a_1 + \frac{3a_3}{4}\right)T_1 + \left(\frac{a_2 + a_4}{2}\right)T_2 + \frac{a_3}{4}T_3 + \frac{a_4}{8}T_4.$$

Since $|T_n(x)| = |\cos(n \arccos x)| < 1$, and with the last member estimated under the condition (2), then we have to obtain:

$$\left|\frac{a_4}{8}\right| < 0,05$$

Consequently, for all fourth degree polynomials to which , the following is valid:

$$f(x) = \left(a_0 + \frac{a_2}{2} + \frac{3a_4}{8}\right)T_0 + \left(a_1 + \frac{3a_3}{4}\right)T_1 + \left(\frac{a_2 + a_4}{2}\right)T_2 + \frac{a_3}{4}T_3 =$$

$$= \left(a_0 + \frac{a_2}{2} + \frac{3a_4}{8}\right) + \left(a_1 + \frac{3a_3}{4}\right)x + \left(\frac{a_2 + a_4}{2}\right)(2x^2 - 1) + \frac{a_3}{4}(4x^3 - 3x) =$$

$$= \left(a_0 - \frac{a_4}{8}\right) + a_1x + (a_2 + a_4)x^2 + a_3x^3$$

Example: A capital investment intended for the production of series A based on the 15 quarterly values of production (variable x in comb.) is being analysed. Variable y (in 000 kn) denotes total investments. Here we have:

x_i - quarterly value of production in kunas for $i = 1, 2, \dots, 15$
 y_i - total investment of dedicated production of series A, expressed in 000 kn.

The regression fourth degree polynomial is chosen as the total investments model. Using the

method of least squares, parameters based on empirical values are estimated, and estimation standard errors are calculated. The regression equation is the following:

$$\hat{y} = 2434 + 85,7x - 0,03x^2 + 0,00004x^3 - 0,000002x^4$$

$$R^2 = 0,999.$$

Since $|a_4| = 0,000002 < 0,4$, the assumption is fulfilled, hence the obtained regression polynomial can be approximated by a third degree polynomial. Consequently, we have:

$$\begin{aligned} P_3(x) &= 2434 + 85,7x - (0,03 + 0,000002)x^2 + 0,00004x^3 = \\ &= 2434 + 85,7x - 0,030002x^2 + 0,00004x^3 \end{aligned}$$

We have obtained the third degree polynomial which approximates the regression fourth degree polynomial with prespecified level of significance and its coefficients still remaining within the calculated confidence intervals.

Conclusion

This paper presents the regression model analysis and variance analysis. When finding the optimal degree of a polynomial, and testing its coefficients in the variance analysis, it is very important to find a polynomial of the least degree possible as the regression polynomial of a higher degree may result in a higher order matrix, which would make matrix multiplication and the search for the inverse matrix more complicated. The results of special functions, i.e. Chebyshev polynomials, have been used in the paper. These polynomials approximate the polynomial obtained by the regression analysis to a lower degree polynomial with required accuracy, keeping the polynomial coefficient within the estimated intervals at the same time. The advantages of the result achieved are obvious because the method enables an easier and faster regression coefficient calculation. In addition, it

avoids unnecessary difficulties in methodology of searching regression values.

Literature

1. Boyd, J. P.: *Chebyshev and Fourier Spectral Methods*, University of Michigan, Dover Publications, 2000.
2. Demeyer, J.: *Recursively enumerable sets of polynomials over a finite field*, J. Algebra 310, 2007.
3. Ditlevsen, O., and Madsen, M. O., *Structural reliability methods*, Technical University of Denmark, Lyngby, (2nd ed.), 2003.
4. Eisentr, K.: *Hilbert's tenth problem for algebraic function fields of characteristic 2*, Pacific J. Math. 210, 2003.
5. Sarapa, N., *Teorija vjerojatnosti*, Školska knjiga, Zagreb, 2002.
6. Tomašević, M., *Matematički i vektorski račun*, Pomorski fakultet Sveučilišta u Splitu, 2007.
7. Tomašević, M., Ristov, P., and Stanivuk, T., *Metodologija znanstvenoistraživačkog rada, Statističke metode u istraživanju*, Pomorski fakultet Sveučilišta u Splitu, 2007.
8. Vasilyev, N. and Zelevinsky, A. *A Chebyshev Polyplayground: Recurrence Relations Applied to a Famous Set of Formulas*. Quantum 10, 20-26, Sept/Oct.1999.
9. Vujanović, N., *Teorija pouzdanosti tehničkih sistema*, Vojnoizdavački novinarski centar, Beograd 1990.
10. <http://mathworld.wolfram.com/ChebyshevPolynomialoftheFirstKind.html> (accessed 2009-07-08)