



Multicollinearity

Aylin Alin*

Multicollinearity refers to the linear relation among two or more variables. It is a data problem which may cause serious difficulty with the reliability of the estimates of the model parameters. In this article, multicollinearity among the explanatory variables in the multiple linear regression model is considered. Its effects on the linear regression model and some multicollinearity diagnostics for this model are presented. © 2010 John Wiley & Sons, Inc. *WIREs Comp Stat* 2010 2 370–374

Multicollinearity refers to the linear relationship among two or more variables, which also means lack of orthogonality among them. This relation is also called *collinearity* or *ill conditioning* by some authors like Belsley¹ and Chatterjee and Hadi.² In more technical terms, *multicollinearity* occurs if k vectors lie in a subspace of dimension less than k . This is the definition of *exact multicollinearity* or *exact linear dependence*. It is not necessary for multicollinearity to be exact in order to cause a problem. It is enough to have k variables nearly dependent, which occurs if the angle between one variable and its orthogonal projections onto others is small. Through this article, multicollinearity will mean near dependence. Multicollinearity is a condition of deficient data, which frequently is encountered in observational studies in which the investigator does not interfere with the study. In many studies, multicollinearity has been confused with correlation. Correlation is the linear relationship between only two variables, whereas multicollinearity can exist between two variables or between one variable and linear combination of the others. Thus, correlation is a special case of multicollinearity. High correlation implies multicollinearity, but the converse is not true. One can have multicollinearity among explanatory variables, but still not have high correlation between pairs of these variables.

Multicollinearity creates difficulties when one builds a regression model between response variable y and explanatory variable X . Through this article, the focus will be on the multicollinearity problem among the columns of X and the multiple linear regression

between y and X . The multiple linear regression model between y and X is given in Eq. (1).

$$y = X\beta + \varepsilon \quad (1)$$

where y is the $n \times 1$ vector of response variable, X is the matrix of explanatory variables with dimensions $n \times (k + 1)$ with the first column consisting of the vector of ones, β is the $(k + 1) \times 1$ vector for regression coefficients, and ε is the $n \times 1$ vector of errors, which is assumed to have mean 0 and variance covariance matrix $\sigma^2 I$. If the explanatory variables are centered, the model given in Eq. (1) does not include intercept and then there is no need to include the vector of ones in X matrix. In this case, the X matrix will have k columns whereas β will have k rows. The least squares estimates of β is obtained as follows:

$$b = (X'X)^{-1}X'Y \quad (2)$$

EFFECTS OF MULTICOLLINEARITY

If there is multicollinearity, the signs of b_i 's for $i = 1, 2, \dots, k$ may be wrong, i.e., different from the signs of correlation between the corresponding explanatory variable and response variable. However, the signs of the regression coefficients can be 'right' even when there is a high degree of multicollinearity and 'wrong' when there is essentially no multicollinearity (Ref 3, p. 131). The most serious effect of multicollinearity is that b_i 's will have large standard errors ($\sigma_{b_i}^2$), i.e., large sampling variability. This makes these coefficients unreliable and, therefore, decreases their precision. The more the multicollinearity is exact, the less reliable are the estimates. In addition to having unreliable coefficients, which may vary from one sample to other, the inflated variances of coefficient estimates harm hypothesis testing, estimation, and forecasting. The regression coefficient b_i is interpreted

Additional Supporting Information may be found in the online version of this article.

*Correspondence to: aylin.alin@deu.edu.tr

Department of Statistics, Dokuz Eylul University, Izmir, Turkey

DOI: 10.1002/wics.84

as the change in the expected value of y when x_i is increased by one unit under the condition that all remaining explanatory variables are held constant. The b_i is called *partial regression coefficient*. This interpretation is not applicable in case of multicollinearity, because it may not be possible in practice to vary one explanatory variable and hold other constant. Therefore, multicollinearity makes it hard to measure the marginal effect of explanatory variables.

In multiple linear regression, extra (sequential) sum of squares (SSR) measures the marginal contribution of an explanatory variable in reducing the error sum of squares (SSE). When there is multicollinearity among explanatory variables, the marginal contribution of any one variable in reducing the SSE depends on which other variables are already in the regression model. For instance, consider the data given by Wold.⁴ Data consist of seven explanatory variables. When there is only x_1 in the model, $SSR(x_1) = 56.964$, but when x_2 and x_3 are in the model, the marginal contribution of x_1 is $SSR(x_1|x_2, x_3) = 1.308$, which is very small when compared with that of $SSR(x_1)$. The explanation for this is the multicollinearity problem among these variables, which will be shown numerically in the following section. The marginal contribution of x_1 is small when x_2 and x_3 are already in the model, because x_2 and x_3 contain much of the same information as x_1 . Including a variable having multicollinearity with others already in the model contributes nothing except decreasing the degrees of freedom. Combining large sampling variability of the coefficient estimate and small marginal contribution to the reduction of SSE, the partial regression coefficients will be insignificant.

Multicollinearity does not affect the model significance given in Eq. (1). Actually, even though all variables in the model are insignificant, we can still have significant model. For the data mentioned above, the P -value for the analysis of variance performed for model significance is 0.007 even though the P -values for the significance of partial regression coefficients are 0.584, 0.067, and 0.615 for x_1 , x_2 , and x_3 , respectively. For model significance, we test the null hypothesis $H_0: \beta_1 = \beta_2 = \dots = \beta_k = 0$. Hence, the reduced model does not contain any of the variables and it does not matter whether there is multicollinearity or not. However, the null hypothesis for the significance of i th partial regression coefficient is $H_0: \beta_i = 0$. The reduced model will include all the variables already in the model except x_i . Unlike the model significance, the multicollinearity will matter this time. However, multicollinearity does not have to affect the precision of the mean response or prediction

of new observations if these inferences are made within the range of observations (see the example in Ref 5, p. 293).

Because multicollinearity is not a modeling error, the problems, which may be caused by multicollinearity, should be investigated after the model has been satisfactorily specified. Belsley¹ (p. 26) emphasizes that some of the problems caused by multicollinearity such as the sensitivity of the regression coefficients to the addition or deletion of a variable or observation or having explanatory variables with low t -statistics may be the result of combination of building the wrong model and having influential data without any presence of multicollinearity.

DIAGNOSING MULTICOLLINEARITY

The effects of multicollinearity mentioned in the previous section can also be used as an indication of multicollinearity, but this would be informal. The mostly used diagnosis is to examine the correlation matrix of explanatory variables. As correlation and collinearity are not same, there can still be multicollinearity even when all correlations are low. Hence, using correlation matrix \mathbf{R} is not enough. Determinant of \mathbf{R} ($\det \mathbf{R}$) is another diagnostic for multicollinearity but it suffers from the same problem with pair-wise correlations. Even though determinant is very close to one, there still might be multicollinearity among the columns of \mathbf{X} . Furthermore, both pair-wise correlations and $\det \mathbf{R}$ do not reveal the number of coexisting relations and their structure.

Variance inflation factor (VIF) given in Eq. (3) is very popular as a multicollinearity diagnostic.

$$VIF_i = \frac{1}{1 - R_i^2} \text{ for } i = 1, 2, \dots, k \quad (3)$$

R_i^2 is the coefficient of multiple determination of x_i on the remaining explanatory variables. Diagonal elements of $(\mathbf{X}'\mathbf{X})^{-1}$ are equal to VIF values when \mathbf{X} is standardized. VIF measures the increase in $\sigma_{b_i}^2$ because of multicollinearity relative to the variance that would result if there was no multicollinearity. The larger the VIF, the more $\sigma_{b_i}^2$ is inflated. A large value of VIF_i indicates the involvement of x_i in at least some linear dependency, but does not reveal which x_j . The threshold value to deviate small from large is generally taken as 10. The mean of VIF values is also used as an index for multicollinearity. The mean considerably larger than one is an indication of multicollinearity. For the data including seven explanatory variables given by Wold⁴, the VIF values are 307.1, 24.1, 25.7,

TABLE 1 | Eigenvalues, Eigenvectors and Condition Indexes for the Data Given in Ref 15

λ	0.027	0.133	0.421	0.555	1.512	40.047	83.305
	−0.659	0.128	0.384	0.397	−0.088	−0.210	0.439
	−0.012	−0.001	0.234	−0.805	0.217	−0.276	0.417
	−0.052	0.052	0.653	−0.252	−0.516	0.241	−0.425
v	0.102	−0.798	0.182	0.089	0.007	0.432	0.354
	−0.125	0.197	−0.477	−0.216	−0.640	0.325	0.396
	0.561	−0.112	0.114	0.208	−0.472	−0.593	0.207
	0.472	0.542	0.313	0.182	0.215	0.422	0.358
η	55.947	25.032	14.071	12.248	7.423	1.442	1.000

TABLE 2 | Variance Decomposition Proportions for the Data Given in Ref 15

η	$\sigma_{b_1}^2$	$\sigma_{b_2}^2$	$\sigma_{b_3}^2$	$\sigma_{b_4}^2$	$\sigma_{b_5}^2$	$\sigma_{b_6}^2$	$\sigma_{b_7}^2$
55.947	0.955	0.004	0.072	0.074	0.329	0.970	0.768
25.032	0.007	0.000	0.014	0.907	0.164	0.008	0.202
14.071	0.021	0.097	0.709	0.015	0.305	0.003	0.021
12.248	0.017	0.872	0.080	0.003	0.047	0.006	0.005

95.1, 32, 219.1, and 196.4 for x_1 , x_2 , x_3 , x_4 , x_5 , x_6 , and x_7 , respectively. Very large VIF values are indicators of multicollinearity problem for this data. Moreover, the mean of these values is 128.5, which is much larger than one. Despite its popularity, VIF is still unable to distinguish the number of coexisting linear dependencies and their structure.

Examining the eigenvalues (λ_s) and eigenvectors (v_s) of R or $X'X$ yields another method to diagnose multicollinearity for $s = 1, 2, \dots, k$. This method, also called a *principal component approach*, reveals not only the number of linear dependencies but also the structural relations within these relations. Small λ values indicate the presence of multicollinearity and there will be one small λ value for each linear dependency. The problem is to determine what small is. Chatterjee and Hadi² (p. 243) suggest that if one λ value is much smaller than the others (and near zero), then multicollinearity is present. The elements of v corresponding small λ value help to diagnose the variables involved in the corresponding linear relation. The nonzero elements of v indicate the variables which are in the linear relation, whereas zero elements indicate which are not. However, Belsley¹ (p. 36) argues this criterion. He claims that it is possible for an element of an eigenvector to be arbitrarily small even though the corresponding variable still belongs to the linear relation. Examining the eigenvalues of $X'X$ is equivalent to examining singular values of $X(\mu_s)$ where $\mu_s = (\lambda_s)^{1/2}$ (see Ref 6, pp. 91–103, 334–336

and Ref 1, pp. 41–44 for more information about eigenvalues, eigenvectors, and singular values). The sum of reciprocals of λ_s for $s = 1, 2, \dots, k$ is also used as multicollinearity diagnostic. The rule of thumb is that if this sum is larger than $5k$ then multicollinearity is present.

Condition number given in Eq. (4) is another measure to diagnose the overall multicollinearity. $\kappa \gg 1$ indicates evidence of strong multicollinearity.

$$\kappa = \sqrt{\frac{\lambda_{\max} \text{ of } X'X}{\lambda_{\min} \text{ of } X'X}} = \frac{\mu_{\max} \text{ of } X}{\mu_{\min} \text{ of } X} \quad (4)$$

κ is the largest among the condition indexes $\eta_s = \mu_{\max} \text{ of } X / \mu_s \text{ of } X$ defined by Belsley¹ (p. 55). η_s around 5–10 indicates weak multicollinearity and 30–100 indicates moderate to strong multicollinearity.

It is mentioned that multicollinearity causes inflation at $\sigma_{b_i}^2$. Variance decomposition proportions are the multicollinearity diagnostics developed based on this idea. These proportions are defined as

$$\pi_{si} = \frac{c_{is}}{\sum_{s=1}^k c_{is}}, \text{ where } c_{is} = \frac{v_{is}^2}{\lambda_s} \quad (5)$$

v_{is} is the i th element in v_s . π_{si} is the proportion of VIF of the i th regression coefficient associated with the multicollinearity represented with λ_s . In other

words, it shows the extent to which $\sigma_{b_i}^2$ is inflated by multicollinearity corresponding to small λ_s or equivalently large η_s . As eigenvector approach, variance decomposition proportions enable us to determine which variables are involved in which linear relations.

Table 1 includes the eigenvalues, eigenvectors, and condition indexes for the data given by Wold.⁴ These values have been obtained using variables centered and scaled to have zero mean and unit length to make the calculations easier and the comparisons more meaningful. The λ values and condition indexes suggest that there are four coexisting linear dependencies in these data from weak to stronger. The elements of v given in bold represent the variables involved in these dependencies. It is clear that x_1 , x_6 , and x_7 are involved in the strongest linear relation. Variance proportions for the four linear dependencies presented with Table 2 reveal that this dominating linear relation accounts for 96, 97, and 77% of the inflation at $\sigma_{b_1}^2$, $\sigma_{b_6}^2$, and $\sigma_{b_7}^2$, respectively. These are all given in bold in Table 2. The values in the eigenvector corresponding to the second smallest eigenvalue reveal that x_4 and x_7 are included in this relation and the variance proportions show that 91 and 20% of the inflation at $\sigma_{b_4}^2$ and $\sigma_{b_7}^2$, respectively are accounted by this relation.

There are other diagnostic procedures, which are less popular such as the ones proposed by Farrar and Glauber⁷ and Haitovsky.⁸ These procedures are based on statistical hypothesis testing. However, they have been criticized by some authors, e.g., Kumar,⁹ Wichers¹⁰ and Parker and Smith.¹¹

It should be noted that multicollinearity may not be harmful. Harmful effects of multicollinearity could be counteracted by a sufficiently small error variance σ^2 . The diagnostics given in this text only include information in X . We need different diagnostics incorporating the size of σ^2 . Belsley¹ suggested signal to noise ratios to detect harmful multicollinearity.

Most of the authors suggest that effects of multicollinearity may be removed by having more data or introducing prior information. Harvey¹² proposes that multicollinearity may be reduced by imposing appropriate restrictions, though it cannot, in general, be eliminated. Soper¹³ proposes a technique to reduce multicollinearity in linear regression, which is also discussed by Kennedy.¹⁴ There are other methods alternative to linear regression when multicollinearity is present in X such as ridge regression, principal component regression, and partial least squares regression. Even though these methods give biased estimates, these estimates have smaller variance compared with ordinary least squares estimates. See Refs 15–18 for more information about these methods.

CONCLUSION

Multicollinearity problem among the explanatory variables, its effects and diagnostics to determine multicollinearity have been presented. As in most of the text books and articles, multicollinearity in the linear regression has been studied in this article. However, it is also a problem for other generalized linear models such as the logistic regression, which is also a very popular technique, especially for medical studies. See Ryan³ and Hosmer and Lemeshow¹⁹ for the detection of multicollinearity in logistic regression.

The multicollinearity diagnostics are available in most of the statistical software. The correlation matrix and VIF values may be obtained easily using Minitab, SAS, and SPSS. SAS also provides eigenvalues, eigenvectors, and variance decomposition proportions. Moreover, MATLAB code for a function calculating eigenvalues, eigenvectors, condition indexes, and variance decomposition proportions will be provided in the electronic version of this text along with options for using original X matrix, scaled but not centered X matrix, or both centered and scaled X matrix.

REFERENCES

1. Belsley DA. *Conditioning diagnostics collinearity and weak data in regression*. New York: John Wiley & Sons; 1991.
2. Chatterjee S, Hadi AS. *Regression Analysis by Example*. 4th ed. New York: John Wiley & Sons; 2006.
3. Ryan TP. *Modern Regression Methods*. New York: John Wiley & Sons; 1997.
4. Wold S, Sjöström M, Eriksson L. PLS-regression: a basic tool of chemometrics. *Chemom Intell Lab Syst* 2001, 58:109–130.
5. Neter J, Kutner MH, Nachtsheim CJ, Wasserman W. *Applied Linear Statistical Models*. 4th ed. Chicago: Irwin; 1990.
6. Seber GAF. *A Matrix Handbook for Statisticians*. New York: John Wiley & Sons; 2008.
7. Farrar DE, Glauber RR. Multicollinearity in regression analysis: the problem revisited. *Rev Econ Stat* 1967, 49(1):92–107.
8. Haitovsky Y. Multicollinearity in regression analysis: comment. *Rev Econ Stat* 1969, 51(4):486–489.

9. Kumar TK. Multicollinearity in regression analysis. *Rev Econ Stat* 1975, 57(3):365–366.
10. Wichers CR. The detection of multicollinearity: a comment. *Rev Econ Stat* 1975, 57(3):366–368.
11. Parker RN, Smith MD. High correlations or multicollinearity, and what to do about either: reply to light. *Soc Forces* 1984, 62(3):804–807.
12. Harvey AC. Some comments on multicollinearity in regression. *Appl Stat* 1977, 26(2):188–191.
13. Soper JC. Second generation research in economic education: problems of specification and interdependence. *J Econ Educ* 1976, 8(2):40–48.
14. Kennedy PE. Eliminating problems caused by multicollinearity: a warning. *J Econ Educ* 1982, 13(1):62–64.
15. Geladi P, Kowalski BR. Partial least squares regression: a tutorial. *Anal Chim Acta* 1986, 185:1–17.
16. Höskuldsson A. PLS regression methods. *J Chemom* 1988, 2:211–228.
17. Frank IE, Friedman JH. A statistical view of some chemometrics regression tools. *Technometrics* 1993, 35(2):109–148.
18. McDonald GC. Ridge regression. *WIREs Comput Stat* 2009, 1:93–100.
19. Hosmer DW, Lemeshow S. *Applied Logistic Regression*. New York: John Wiley & Sons; 1989.