



## **Ridge Estimates of Regression Coefficients for Soil Moisture Retention of Iraqi Soils**

Khasraw A. Rashid<sup>1</sup> & Hanaw A. Amin<sup>2</sup>

*1 Department of Soil and Water Science, Faculty of Agricultural Science, University of Sulaimani*

*2 Faculty of Science and Science Education, University of Sulaimani.*

*E- Mail: [khasraw.rashid@univsul.edu.iq](mailto:khasraw.rashid@univsul.edu.iq)*

### **Article info**

Original:  
08 October 2015  
Revised:  
31 January 2016  
Accepted:  
25 February 2016  
Published online:  
20 September 2016

**Key Words:** Multiple Regression Model, Ridge Regression Model, Collinearly, SCSSNI.

### **Abstract**

Statistical literature has several methods for coping with multicollinearity problem. Ridge regression “RR” is compared with multiple linear regression for studying “Suction Characteristics of Subgrade Soils from North Iraq- SCSSNI”. Multiple linear regression “MLR” of “SCSSNI” data usually encounters a collinearly problem, which adversely affects long term prediction performance. The collinearly problem can be eliminated or greatly improved by using ridge regression, which is a biased estimation method with potentially smaller mean square error “MSE” as an alternative to ordinary least square “OLS”. In this study, ridge regression (a biased estimation method has been evaluated with a constant bias ( $k$ )) and the prediction performance was compared with that of ordinary least square “OLS” based multiple linear regression “MLR”. The bias constant of ( $k = 0.100000$ ) was selected by examining the ridge trace. At this point, the estimated coefficients are stable and their variance inflation factors “VIFs” become smaller. To evaluate the robustness of each model the standard error of prediction “SEPs” has been compared, the prediction of original values using MLR model shows slightly better results comparing to that of ridge regression model, which is due to an intentional bias is associated in the ridge model.

To compare RR and MLR, the coefficient of determination, “VIFs”, and standard error “SE” of parameters has been studied. If the variance of the ridge estimator  $\hat{B}_R$  could be tremendously reduced, the mean square error tends to be smaller than the OLS. The prediction results of a ridge model showed more stable prediction performance as compared to that of MLR, by removing or decreasing the collinearly problem.

### **Introduction:**

Soil suction, in a very general manner, may be looked upon as the “Capacity Potential” of a soil, since water is held in a soil (above the water table) by surface tension and by adsorption force. It has been defined [1] in more precise terms as “the work required to move a unit mass of water against capillary forces in a column of soil, from a free water surface to a given point above this surface”. [2]

Curves showing the relation between moisture content and soil moisture tension have been referred to as “sorption curves” [3], “characteristic curves” [4], “retention curves” [5].

Multiple linear regression is one of the most popular calibration methods for data analyzing, in comparison to other calibration methods, is simple, easy to understand, and possibly to clearly rationalize the relation between SCSSNI (suction characteristics of soil from north of Iraq ) and prediction variables (Clay= $X_1$  , Silt= $X_2$ , Sand= $X_3$ , Organic Matter= $X_4$ , Calcium Carbonate= $X_5$ , Cation Exchange Capacity= $X_6$  , Saturation Percentage= $X_7$ ) . There are several methods can be used such as forward, backward and stepwise regression for estimating the statistical parameters. In general, variables in SCSSNI are highly correlated each other (which is referred to as a multicollinearity problem)[6]. Multicollinearity, or collinearity, is the existence of near-linear relationships among the independent variables [7]. In the presence of multicollinearity, estimates of least square methods including MLR are unstable and tend to poor prediction. Multicollinearity can create inaccurate estimates of the regression coefficients, inflate the standard errors of the regression coefficients, deflate the partial t-tests for the regression coefficients, give false, no significant, p-values, and degrade the predictability of the model. For dealing with multicollinearity, it needs to be able to identify its source. The source of the multicollinearity impacts the analysis, the corrections, and the interpretation of the linear model.

There are five sources: (data collection, physical constraints of the linear model or population, over-defined model, model choice or specification, and Outliers) [7, 8]. It is known that biased estimation methods gives considerably better prediction than (OLS) when the data are noisy or the prediction are highly collinear[9]. Ridge regression a (biased estimation method) has been evaluated and the prediction performance was compared with that OLS based MLR. To compare ridge regression and MLR. MLR model was initially developed using OLS methods for estimating regression parameters, and then RR model is used to predict regression parameters.

### **Materials and Methods:**

This study was confined only to the northern part of Iraq. Figure (1) Shows 64 locations from where the various representative types of soil samples were taken for laboratory analysis [6].

The technique used, for measuring soil moisture at different pF, value determined according to the procedure suggested by [6].

The soil samples were collected from varies locations as outlined above were subjected to the following laboratory tests, after air- drying, quartering and passing through a BS 2mm size sieve:

- i) Particle size distribution: According to the method described in[2]
- ii) Soil moisture at different pF: (pF = log column of water, cm) values determined according to the procedure suggested by [5, 10].
- iii) Calcium carbonate content was determined according to the method described in[11].
- iv) Organic matter content units: determined according to the procedure described in[12].
- v) Cation Exchange Capacity(CEC): was determined according to the method given by [13].

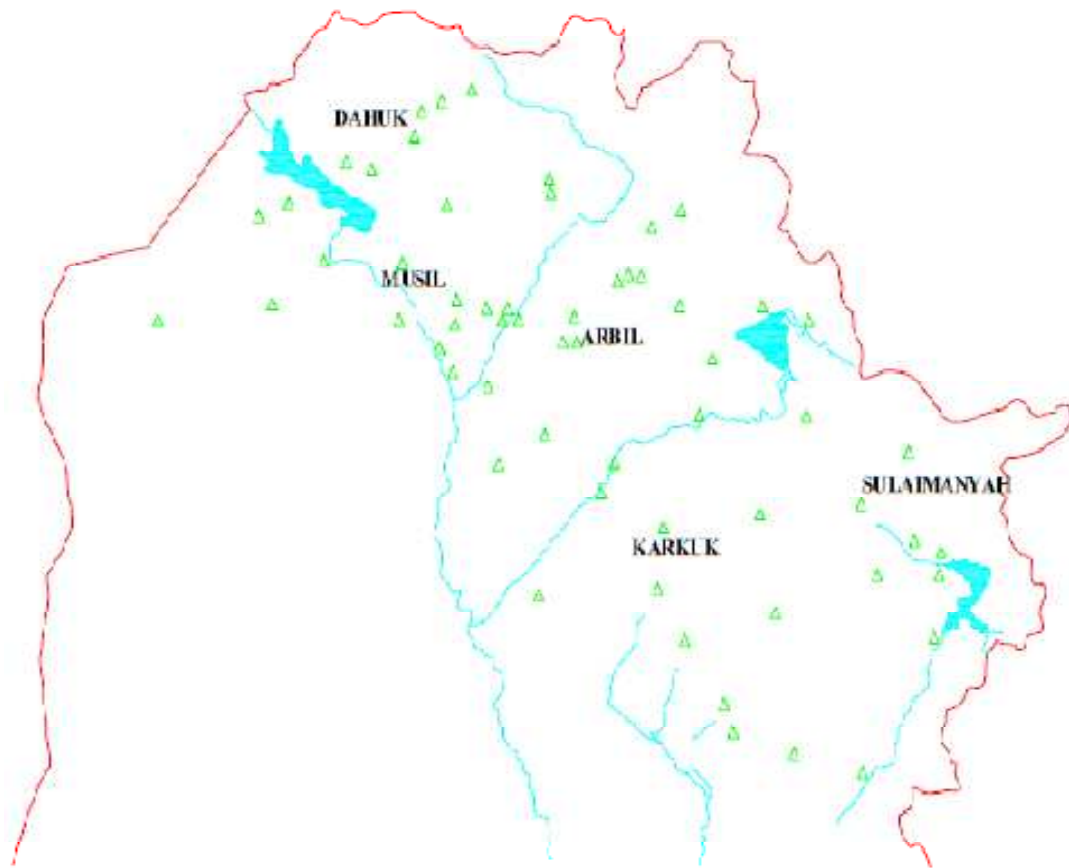


Figure- 1: Shows locations from where the soil samples were taken

A typical pF moisture content curve can be generalized as in Fig. (2a,2b) with parameters (A) and (B) representing the moisture contents at pF = 0.0 and pF = 4.0 respectively. It may be observed that higher values of (A) are associated in general with higher clay contents, similarly higher values of (B) are associated with lower values of sand content. Further higher values of (A) and clay content are associated with higher values of CEC [Cation Exchange Capacity]. It is, therefore, obvious that the parameters (A, B) and (A-B) can be used from the pF moisture content plots to indicate type of soil.

Taking a cue from the pF moisture content plots where in A and B parameters indicate the type of soil, a linear regression analysis was carried out. In this analysis, the following are the dependent variable [6]:

$Y$  = Moisture content at pF = 0.0 is a dependent variable [6].

The independent variables selected for study are [6]:

$X_1$  = % clay content

$X_2$  = % silt content

$X_3$  = % sand content

$X_4$  = % organic matter content

$X_5$  = % calcium carbonate content

$X_6$  = % cation exchange capacity

$X_7$  = % saturation

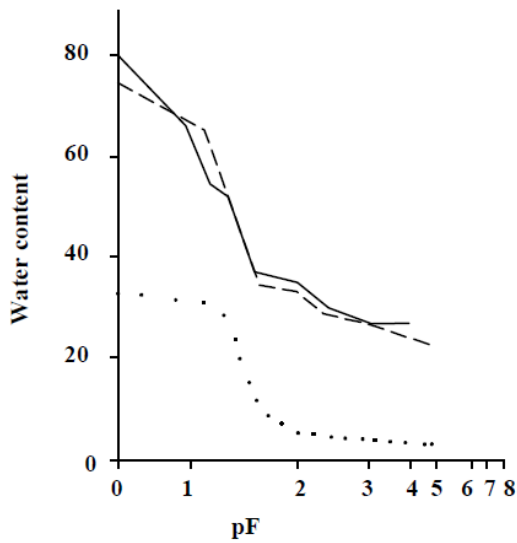


Figure- 2a: pF- moisture content curve for different texture soils

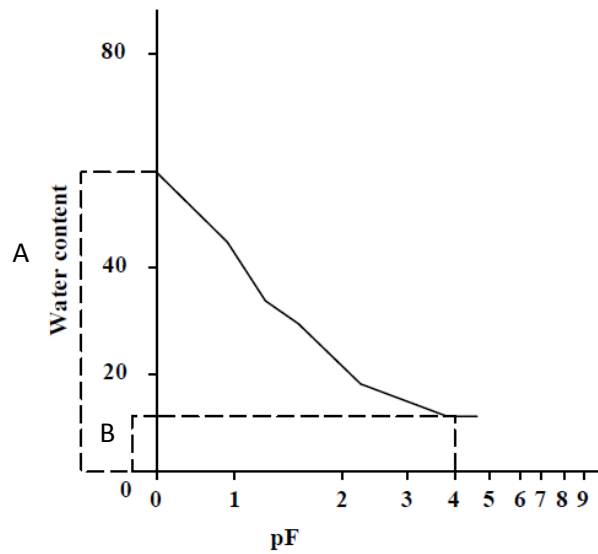


Figure- 2b: shows typical pF- moisture content curve

**Multiple Linear Regression Analysis:**

First we can explain the concept of regression analysis, is a collection of statistical techniques for modeling and investigating the relationship between a response variables of interest (Y) and a set of repressor or prediction variables( $X_1, X_2, \dots, X_n$ ). In the simple linear regression analysis one a relation of the type:[14]

$$Y_i = \beta_0 + \beta_1 X_{i1} + \epsilon_i \tag{1}$$

Where  $i = 1, 2, \dots, n$  and  $\epsilon_i \sim (0, \sigma^2)$  are identically independent containing one explanatory variable. Applications of regression are numerous and occur in almost every applied field including engineers and chemical, physical, science, life and biological science, the social sciences, management and economics. A very important type of regression is multiple linear regression models.

$$Y_{ij} = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_p X_{ip} + \epsilon_{ij} \tag{2}$$

In the response is a linear function of the unknown model parameters or regression coefficients  $\beta_0, \beta_1, \beta_2, \dots, \beta_p$ . Linear regression models are widely used as empirical models to approximate some more complex and usually unknown functional relationship between the response and the regresses variables we can thus write the regression model in the matrix model form as:[9, 14, 15]

$$\underline{Y} = X\underline{\beta} + \underline{\epsilon} \tag{3}$$

Where  $\underline{Y}$  is  $(n \times 1)$  vector of the variable to be explained with  $(n)$  number of observations and  $X$  is a  $(n \times (p + 1))$  design matrix of predictors.  $\underline{\beta}$  is a  $(p + 1)$  vector of regression coefficients  $\beta_j, j = 1, \dots, p$  which needs to be estimated for  $(p)$  independent variables and  $(\underline{\epsilon})$  is a  $(n \times 1)$  vector of normally distributed random errors, with  $\epsilon \sim N(0, \sigma^2)$  identically independent distributed.[14]

The MLR is considered with an intercept having a single response variable(Y)and explanatory variables( $X_1, X_2, \dots, X_p$ ). The estimator vector of regression coefficients by OLS is given by:[9, 14]

$$\hat{\beta}_{OLS} = (X'X)^{-1}X'Y \tag{4}$$

Where  $(X)$  is matrix consisting of the sample data of variables, if  $X_i$ 's are highly correlated, the determinant of  $(X'X)$  is near (zero) and so  $(X'X)$  becomes near- singular.[9] Therefore, in the presence of multicollinearity, the OLS estimates could become very unstable due to large variance of the estimates, which leads to poor prediction.[9, 15]

A theory of OLS was developed by researchers for estimating parameters in a general linear regression model. As it is known the OLS methods is one of the most common procedures that are used in statistics which gives unbiasedness and minimum variance among other estimators, when the basic familiar assumptions are met, if one or more of these assumptions not realized it will lead to invalid conclusions. Some one of these problems concerning linear methods is (heteroscedasticity, autocorrelarity, and multicollinearity).[16]

This method provides the best linear unbiased estimator- BLUE, which is equivalent to the maximum likelihood estimator- MLE for estimation of the underlying normal regression relationship. [16] Also the principle components analysis is one of the simple multivariate methods and searches for a few linear combinations which can be used to summarize the data, losing in the process as little information as possible, it, therefore, can solve the multicollinearity problem.

**Ridge Regression Analysis:**

Ridge regression is one of the several methods to overcome the multicollinearity problem by modifying the OLS to allow a small bias via a constant (k) in the parameter estimate. Ridge regression is like least squares but shrinks the estimated coefficients towards zero. Hoerl and Kennard in (1970) suggested another biased RR estimator with potentially smaller MSE as an alternative to OLS, but this method led to biased estimator.[9, 12, 14, 15, 17]

$$\hat{\beta}_R = (X'X + kI)^{-1}X'Y \dots \dots \dots (5)$$

Where (I) is the identity matrix for the model with the (i) possible predictors ( $X_1, X_2, \dots, X_p$ ). When an estimator has only a small bias and is more precise than an unbiased estimator, it will be closer to the true parameter's value. If the variance of ridge estimator ( $\hat{\beta}_R$ ) could be tremendously reduced, the mean square error tends to be smaller than OLS.[5] As with least squares, ridge regression seeks coefficient estimates that fit the data well, by making the RSS (regression sum of squares) small.[18]

In this situation the point estimate ( $\hat{\beta}_R$ ) becomes more stable, and the confidence interval of ( $\hat{\beta}_R$ ) is narrower. In ridge regression ( $X'X$ 's) are recommended to be transformed by the correlation form, which makes the diagonal element of ( $X'X$ ) equal (1) and the off diagonal element represents the correlation coefficient of the two variables. Since the values of all elements are of the same order of magnitude, this would control round-off errors in inverting ( $X'X$ ) to obtain the ridge estimator ( $\hat{\beta}_R$ ).[8, 9]

This equation can have a unique solution even when ( $X'X$ ) is singular. Thus one application of ridge regression is to produce regression estimates for singular design matrices.[8, 10]

One of the most important in ridge regression is the bias constant (k), where (k) is a positive number. In applications, the interesting values of (k) usually lie in the range (0, 1)[9, 15]. The bias constant in ridge regression is used to reduce the variance of estimates of regression coefficients that are due to multicollinearity. Several methods have been proposed for determining the optimal value of (k).[2, 9, 10, 15] A common strategy is to determine the smallest (k) that makes stable coefficients in the ridge trace with the lowest values of VIF the ridge trace is the plot of VIFs versus (k) values. [If too large, the resulting analytical performance will be degraded by applying a large bias, even though the collinearity problem can be solved][9].

To choose an appropriate value of (k) these three points can be looking for:[17]

- Getting the variance inflation factors (VIF) close to (1).
- Estimated coefficients should be stable.
- Looking for only modest change in ( $R^2$ ) or ( $\sigma^2$ ).

The variance inflation factors (VIF) were obtained to examine the degree of multicollinearity. The VIF for ( $i^{th}$ ) regression coefficient is computed as follows:

$$VIF_i = \frac{1}{1-R_i^2} \dots \dots \dots (6)$$

We see that as the R-squared in the denominator gets closer and closer to one, the variance (and thus  $VIF_i$ ) will get larger and larger. The rule of thumb cut-off value for VIF is 10. Equation (7) shows that the variance of ( $i^{th}$ ) regression coefficient is inflated proportional to  $VIF_i$ : [9, 17]

$$Var.(\hat{\beta}) = \frac{\sigma^2}{\sum_j x_{ij}^2} (VIF_i) \dots \dots \dots (7)$$

Where  $x_{ij}^2$  is the ( $j^{th}$ ) centered sample value of the ( $i^{th}$ ) independent variable and  $\sigma^2$  is the variance of error terms in the MLR model. As a rule of thumb, if  $VIF$  exceeds (30 or 10), it is an indication that the associated coefficients are poorly estimated because of multicollinearity. [9]

To evaluate the robustness of each model MLR and RR the resulting SEPs (Standard Error of Prediction) are compared: [9]

$$SEP = \sqrt{\sum_{i=1}^n (\hat{y}_i - y_i)^2 / n} \dots \dots \dots (8)$$

The assumptions are the same as those used in regular multiple regression: linearity, constant variance (no outliers), and independence. Since ridge regression does not provide confidence limits, normality need not be assumed [7].

**Results and Discussion:**

Since these quantitative variables are not independent, the correlated makes calculations not easier in multiple linear regression analysis type. [15]

To illustrate the ideas to this point, as well as to suggest how regression may have useful applications in soil samples, according to the data and to compare the performance of MLR and RR these baseline were summarized under testing the null hypothesis  $H_0: \beta_j = 0$  against the general alternative  $H_1: \beta_j \neq 0$  for both model equation (4) and (5) by using both program package SPSS and NCSS97 the following results was obtained:

$$\hat{Y}_i = -37.957 + 0.451X_1 + 0.509X_2 + 0.478X_3 + 0.777X_4 + 0.010X_5 + 0.912X_6 + 0.354X_7$$

Table- 1: Analysis of variance using multiple linear regression for soil water content at pF=0

		ANOVA <sup>b</sup>				
Model		Sum of Squares	df	Mean Square	F	Sig.
1	Regression	6019.432	7	859.919	33.445	.000 <sup>a</sup>
	Residual	1439.840	56	25.711		
	Total	7459.271	63			
	R- square	.807				
	Adjusted R-square	0.783				
	$R_{Y,X1,\dots,X7}$	0.898				

a. Predictors: (Constant), X7, X6, X2, X5, X4, X1, X3

b. Dependent Variable:  $Y_i$

The results in the Table (1) are by testing the hypothesis of  $H_0: \beta_j = 0$ . Since (Sig. = 0.000) the null hypothesis is rejected. It can be conclude that there is statistically relation between soilmoisture tension components ofsoils samplesat ( $\alpha = 0.05$ ) levels of significance. It may be concluded that there is a positive trend between components of soil samples as shown in Fig. (3).

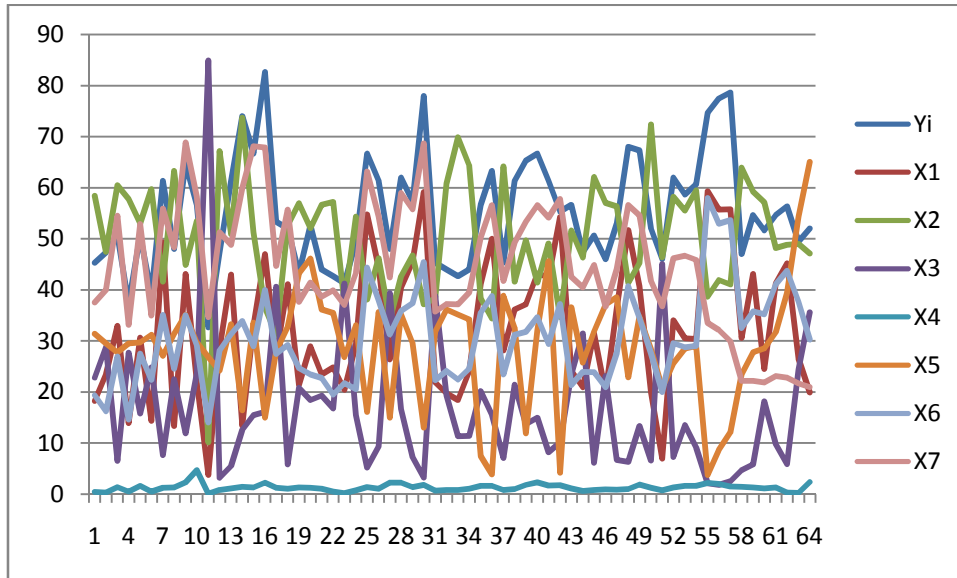


Figure- 3: trend of soil samples contents individuals

The values of coefficient of determination ( $R^2 = 81\%$ ) and adjusted  $R^2 = 78\%$  from the regression of  $X_i$  on the other independent variables. This means that the estimated regression equation (2) of soil samples can be explained by the contents at 81% and adjusted  $R^2 = 78\%$  level percentage.

A number of linear regression equations were developed relating each of the dependent variable  $Y_i$  with each of the independent variables. The expressions which yielded the best statistical fit are:

$$\hat{Y}_i = 27.148 + 0.249X_1 + 0.653X_6$$

The values of coefficient of determination ( $R^2 = 0.674\%$ ) and adjusted  $R^2 = 66.4$  from the regression of  $X_i$  on the other independent variables. This means that the estimated regression equation (2) of soil samples can be explained by the contents at 67% and adjusted  $R^2 = 66\%$  level percentage. As noted the independent variables that yielded the best statistical fit are the % clay content and the cation exchange capacity of soil relating the dependent variable of moisture content  $pF = 0.0$ . The strong positive multiple correlation coefficient of ( $R_{Y,X_1,X_6} = 0.821$ ) gives a statistically significant correlation between the soil properties of the 64 sample from different locations. But the multiple correlation coefficient ( $R_{Y,X_1,\dots,X_7} = 0.898$ ) may be considered to be fairly significant statistically, considering the wide variation in the soil properties of the 64 samples from different locations [17, 19]

Table (2) shows the multicollinearity problem, considering the variance inflation factors (VIF) and R-square. The rule of thumb cut-off value for VIF is 10. VIFs over 10 indicate collinear variables. Solving backwards, this translates into an R-squared value of 0.90. Hence, whenever the R-squared value between one independent variable and the rest is greater than or equal to 0.90, it will be to face multicollinearity.

Table- 2:Least square multicollinearity test for soil samples contents

Independent variables	R- square vs. other X's	VIF of MLR	Tolerance
X <sub>1</sub>	<b>0.9929</b>	<b>139.9576</b>	<b>0.0071</b>
X <sub>2</sub>	<b>0.9905</b>	<b>104.8018</b>	<b>0.0095</b>
X <sub>3</sub>	<b>0.9937</b>	<b>157.9924</b>	<b>0.0063</b>
X <sub>4</sub>	0.4095	1.6934	0.5905
X <sub>5</sub>	0.3048	1.4384	0.6952
X <sub>6</sub>	0.8118	5.3122	0.1882
X <sub>7</sub>	0.5223	2.0934	0.4777

Eigen values of correlation are summarized in Table 3, since some condition numbers greater than (1000). Multicollinearity is a severe problem since (1228.68) is available for X<sub>7</sub>.

Table-3: Eigen values of correlation matrix

No.	Eigen value	Incremental percent	Cumulative percent	Condition number
1	3.076963	43.96	43.96	1
2	1.502264	21.46	65.42	2.05
3	1.001309	14.3	79.72	3.07
4	0.717198	10.25	89.97	4.29
5	0.587222	8.39	98.36	5.24
6	0.11254	1.61	99.96	27.34
7	0.002504	0.04	100	<b>1228.68</b>

The model of ridge regression analysis equation (5) is:

$$\hat{Y}_i = -4.183 + 0.1446X_1 + 0.180X_2 + 0.146X_3 + 1.022X_4 + 0.0095X_5 + 0.878X_6 + 0.331X_7$$

The ridge trace in Fig. (4) is the plot of VIF versus (k) values. This is the famous ridge trace that is the signature of this technique. It presents the standardized regression coefficients on the vertical axis and various values of k along the horizontal axis.

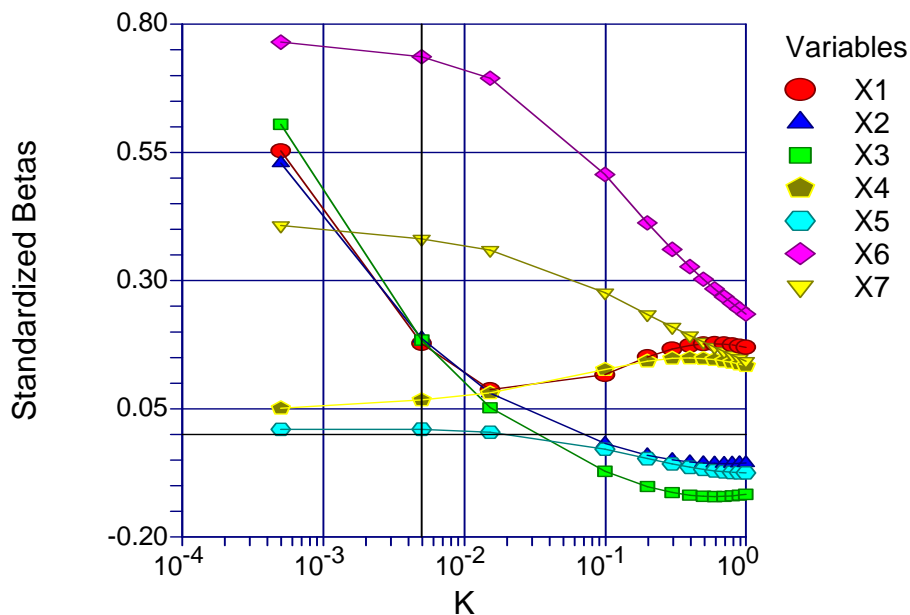


Figure- 4: plot of standardized betas verses various values of (k)

The ridge standard error values in comparison to OLS, when multicollinearity occurs, least squares estimates are unbiased, but their variances are large so they may be far from the true value. By adding a

degree of bias of ( $k = 0.100000$ ) to the regression estimates, the standard errors of soils contents ( $X$ 's) are reduced to almost half as compared to those in MLR.

The results of ridge regression using a bias constant of ( $k = 0.100000$ ) vs OLS are summarized in Table 4. Results, are greatly decreased and more statistically stable model achieved, although regular ridge coefficient and stand'zed ridge coefficients ( $z$ - score) are moderately decreased. So it may be worth sacrificing some bias to achieve a lower variance.

Table-4: Ridge vs. Least Squares Comparison for the values of ( $\beta$ )s

Independent Variable	Regular Ridge Coefficients	Regular L. S. Coefficients	Stand'zed Ridge Coeff's	Stand'zed L.S. Coeff's	RR Standard Error	OLS Standard Error
Intercept	-4.182959	-37.95709				
$X_1$	0.1445993	0.4505841	0.1774	0.5527	0.2017306	0.5661924
$X_2$	0.180119	0.5085251	0.1874	0.5290	0.2016332	0.5778019
$X_3$	0.1455407	0.4781761	0.1840	0.6044	0.2013557	0.5838432
$X_4$	1.021829	0.7772878	0.0674	0.0513	1.124215	1.15773
$X_5$	9.756265E-03	9.505201E-03	0.0102	0.0099	6.737762E-02	6.748981E-02
$X_6$	0.878219	0.9118193	0.7364	0.7646	0.1564915	0.1613807
$X_7$	0.3313024	0.3539966	0.3810	0.4071	6.905199E-02	7.385585E-02
R-Squared	0.8018	0.8070				
Sigma	5.1378	5.0706				

Fig. (5) and Table (5) shows the relation between ( $k$ ) and VIF as increasing ( $k$ ) in the seven- variable model. The VIFs in the three variables ( $x_1, x_2, x_3$ ) are high when ( $k$ ) is near to (0), but they decrease as the value of ( $k$ ) increases. Most noticeably, VIFs are very high at variable ( $x_1, x_2, x_3$ ) and they drop steeply as the bias constant increases.

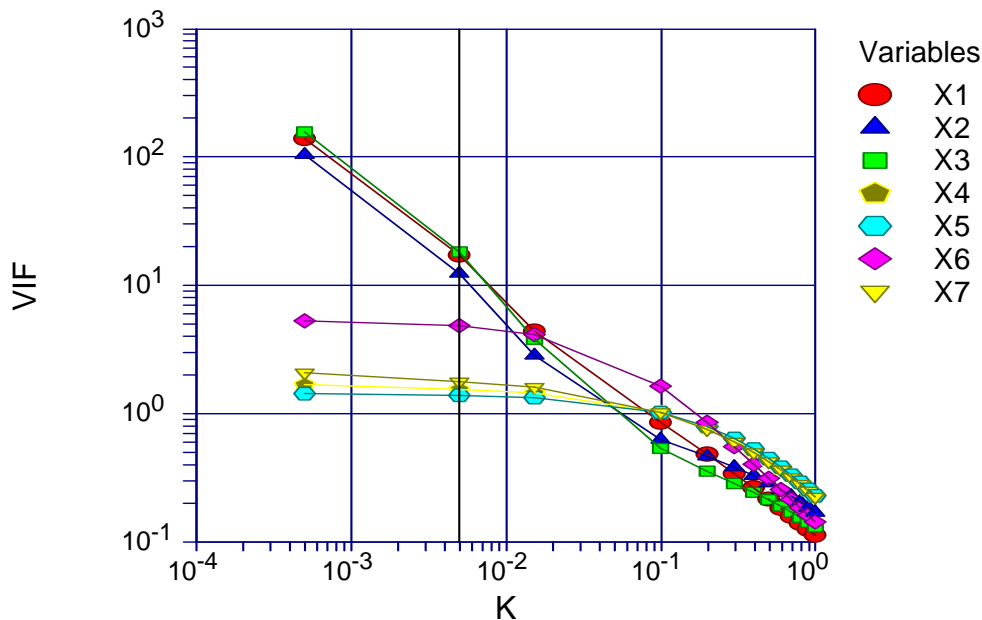


Figure- 5: The relationship between ( $k$ ) bias constant and VIF of each variable as increasing ( $k$ )

In ridge regression, the first step is to standardize the variables (both dependent and independent) by subtracting their means and dividing by their standard deviations as it has been with standard ridge coefficients. As far as standardization is concerned, all ridge regression calculations are based on standardized variables. When the final regression coefficients are displayed, they are adjusted back into their original scale [16, 20]. However, the ridge trace is in a standardized scale. In comparison to MLR results, with a small bias of ( $k = 0.100000$ ), VIF values are greatly decreased and reduced especially at

( $x_1, x_2,$  and  $x_3$ ), and more statistically stable model is achieved, although  $R^2$  is slightly decreased. The result of ridge regression VIF values are summarized in Table (5).

Table (5) the statistical results of soils samples content RR model using the bias constant ( $k = 0.100000$ )

Indep.Var.s	RR Coefficient	SE	Stand. Ridge Coeff.	VIF of RR
Intercept	-4.182959			
X1	0.1445993	0.2017306	0.1774	17.3052
X2	0.180119	0.2016332	0.1874	12.4308
X3	0.1455407	0.2013557	0.184	18.3036
X4	1.021829	1.124215	0.0674	1.5553
X5	9.76E-03	6.74E-02	0.0102	1.3964
X6	0.878219	0.1564915	0.7364	4.8653
X7	0.3313024	6.91E-02	0.381	1.7823

A very important and common strategy is to determine the smallest value of ( $k$ ) for making the stable value of coefficient in the ridge trace with the lowest value of VIF is at ( $k = 0.100000$ ) as explained in Table (6). The fourth row gives an appropriate analysis to give the best statistical fit to RR vs OLS. It can be concluded that for the biased parameter ( $k = 0.100000$ ) the coefficient of determination  $R^2 = 0.7552$ , otherwise R-square is slightly decreased in the other next rows. According to ( $\sigma = 5.7104$ ), but in the other rows is slightly increased, but the average value of  $VIF = 0.9661$  with maximum value of ( $VIF = 1.6457$ ) and it is near to (1). So the best analysis is when the bias parameter ( $k = 0.100000$ ) with parameter estimates seems to make the most sense and it indicates that the associated coefficients of ( $\beta$ )'s are accurate after excluding the multicollinearity problem. Sigma is the square root of the mean squared error. Least squares minimize this value, so we want to select a value of  $k$  that does not stray very much from the least squares value.

Table-6: K Analysis Section

Row	$k$	$R^2$	$\sigma$	$B'B$	Ave VIF	Max VIF
1	0.000000	0.807	5.0706	1.7037	59.0413	157.9924
2	0.005000	0.8018	5.1378	0.7925	8.2341	18.3036
3	0.015312	0.7947	5.2298	0.6343	2.8102	4.4213
<b>4</b>	<b>0.100000</b>	<b>0.7552</b>	<b>5.7104</b>	<b>0.369</b>	<b>0.9661</b>	<b>1.6457</b>
5	0.200000	0.7233	6.0707	0.2822	0.6446	0.8594
6	0.300000	0.6973	6.3499	0.2422	0.4931	0.6447
7	0.400000	0.6744	6.5857	0.2174	0.4001	0.533
8	0.500000	0.6536	6.7927	0.1992	0.3357	0.4492
9	0.600000	0.6344	6.9782	0.1848	0.2881	0.3845
10	0.700000	0.6166	7.1467	0.1727	0.2512	0.3333
11	0.800000	0.5998	7.3009	0.1622	0.2218	0.2921
12	0.900000	0.5841	7.4432	0.1529	0.1978	0.2583
13	1.000000	0.5692	7.5751	0.1446	0.1778	0.2303

Therefore, ridge regression puts further constraints on the parameters,  $\beta_i$ 's, in the linear model. In this case, instead of just minimizing the residual sum of squares also it has a penalty term on the  $\beta$ 's. This penalty term is  $(k - a \text{ pre-chosen constant})$  times the squared norm of the  $\beta$  vector. This means that if the  $\beta_i$ 's take on large values, the optimization function is penalized. It would prefer to take smaller  $\beta$ 's, that are close to zero to drive the penalty term small. [5]

Row (13) in table (7) can be chosen for this purpose at  $(k = 0.10)$ , since this row gives the best line for this case for different values of the biased constant  $(k)$ . Lowest value of VIF is at  $(k = 0.10)$  as explained in table (7) the VIFs values are greatly decreased and a more statistically stable as at  $(k)$  increase, by adding a degree of bias to the regression estimates, ridge regression reduces the standard errors. If the variance of the ridge estimator  $\hat{\beta}_R$  could be tremendously reduced, the mean square error tends to be smaller than the OLS.

Table-7: Variance inflation factor with different values of  $(k)$

row	$k$	$X_1$	$X_2$	$X_3$	$X_4$	$X_5$	$X_6$	$X_7$
1	0	139.9576	104.8018	157.9924	1.6934	1.4384	5.3122	2.0934
2	0.005	17.3052	12.4308	18.3036	1.5553	1.3964	4.8653	1.7823
3	0.01	7.3506	4.9984	7.0674	1.502	1.3682	4.4951	1.6907
4	0.015312	4.4213	2.8567	3.8327	1.4551	1.3406	4.1496	1.6157
5	0.02	3.3071	2.0688	2.6448	1.4178	1.3175	3.8785	1.558
6	0.03	2.2406	1.3582	1.5785	1.3473	1.2716	3.3867	1.4524
7	0.04	1.7557	1.0675	1.1474	1.2861	1.2294	2.9877	1.3642
8	0.05	1.4705	0.9133	0.9227	1.232	1.1901	2.6594	1.289
9	0.06	1.2772	0.8175	0.7862	1.1837	1.1535	2.3858	1.2237
10	0.07	1.1347	0.7514	0.6946	1.1399	1.1191	2.1554	1.1665
11	0.08	1.0237	0.7023	0.6285	1.1	1.0867	1.9594	1.1157
12	0.09	0.9342	0.6638	0.5784	1.0633	1.0561	1.7912	1.0703
<b>13</b>	<b>0.1</b>	<b>0.86</b>	<b>0.6325</b>	<b>0.5389</b>	<b>1.0293</b>	<b>1.027</b>	<b>1.6457</b>	<b>1.0292</b>
14	0.2	0.4862	0.4672	0.3578	0.784	0.7993	0.8594	0.7582
15	0.3	0.3421	0.3856	0.2872	0.6299	0.6447	0.5569	0.605
16	0.4	0.2651	0.3303	0.2447	0.5214	0.533	0.4047	0.5013
17	0.5	0.2169	0.2887	0.2145	0.4405	0.4492	0.3151	0.4253
18	0.6	0.1836	0.2557	0.1911	0.3781	0.3845	0.2567	0.3669
19	0.7	0.1592	0.2287	0.1722	0.3287	0.3333	0.2158	0.3205
20	0.8	0.1404	0.2062	0.1564	0.2888	0.2921	0.1857	0.2829
21	0.9	0.1255	0.1871	0.143	0.256	0.2583	0.1626	0.2519
22	1	0.1134	0.1707	0.1315	0.2287	0.2303	0.1443	0.2259

The resulting (SEP)sof robustness of each model was compared. The prediction of original values using MLR model shows slightly better results ( $SEP^{OLS} = 4.70$ ) comparing to that of ridge regression model ( $SEP^{RR} = 4.76$ ), which is due to an intentional bias is associated in the ridge model. Practical results indicate that RR is not always better than other estimators in terms of SEP. RR is best and depends on the ridge parameter k. For suitable estimates of these parameters, RR estimator might be considered as one of the good estimators using SE.

### Conclusions:

- 1- It is concluded that there is statistically relation between soil moisture retention and component of soil sample at ( $\alpha = 0.05$ ) level of significance.
- 2- It is concluded that there is a positive trend between components of soil samples.
- 3- It is concluded that the RR predictors are with greater stability compared to that from MLR and this depend on ridge parameter( $k = 0.1$ ).
- 4- The results show that the RR has a lower SE than the MLR method for different situations that has been evaluated.
- 5- This study is also concluding that OLS should not be used when the data is collinear since the vector of estimate parameters become too long.
- 6- The RR is preferred since it offers some reduction of the MSE.
- 7- Both linear regression and RR are good method for obtaining more stable parameter estimates.

### References:

- [1] Baver, L. D. Gardner W. H and Gardner W. R. "Soil Physics" 4<sup>th</sup> edition, New York, Wiley and Sons, (1972).
- [2] Kilmer, V. J. and L. T. Alexander. "Methods of Making Mechanical Analysis of Soil". Soil Sci. Vol.68, pp 15 – 24.
- [3] Russell, M. B. "Soil Moisture Sorption Curves for Four Iowa Soils". Soil Sci. Amer Proc. Vol.4, pp 51-54, (1940).
- [4] Child, E. C., "The Use of Soil Moisture Characteristics in Soil Studies", soil science Vol. 50, pp 239 – 252, (1940).
- [5] Richards, L. A. A. "Pressure Membrane Extraction Apparatus for Soil Solution. Vol. 15 p. 377 – 389, (1940).
- [6] Rashid, Khasraw Abdulla. "Suction Characteristics of Subgrade Soil from North Iraq".Dep. Of Civil Engineering, University of Salahaddin, Arbil- Iraq. Engineering and Technology Vol. 12 No. 9, (1993).
- [7] James, Gareth, Witten Daniela, Hastie Trevor, and Tibshirani Robert. "An Introduction to Statistical Learning with Applications in R".Springer New York Heidelberg Dordrecht London, (2013).
- [8] Tibshirani, Robert, James Gareth, Witten Daniela, Hastie Trevor. "An Introduction to Statistical Learning with Applications in R".Springer New York Heidelberg Dordrecht London, (2013).
- [9] Chung, Hoeil, and Jun, Chi- Hyuck. "Determination of Research Octane Number Using NIR Spectral Data and Ridge Regression".Department of Industrial Engineering Pohang University of Science and Technology, Korea, (2000).
- [10] Richards, L. A., "Physical Condition of Water in Soil, A me. Soc. Agron, Inc Madison, Wis. P. 128-151, (1965).
- [11] Richards, L. A.."Diagnosis and Improvement of Saline and Alkaline Soils".Agr.Hanbook No. 60, USDA. US. Government Printing Office Washington DC., (1954).
- [12] Wakley, A. and Black, 1934 cited in L. E. Allison, Organic Carbon, P. 1367 – 1378in black C. A. et al.,

- (Ed.) 1965, Methods of Soil Analysis part 2, Agron. No. 9.
- [13] Chapman, H. D. and Partt, P. E. "Method of Analysis for Soils, Plant and Water". Univ. California Press. California (1961).
- [14] Cule, Erika, Vineis Paolo and De Iorio Maria. "Significance Testing in Ridge Regression for Genetic Data". Cule et al. BMC Bioinformatics, (2011).
- [15] Batah, Feras Sh. M. and Damodar, Sharad Gore. "Ridge regression estimator: Combining Unbiased and Ordinary Ridge Regression Methods for Estimations". Surveys in Mathematics and its Applications. ISSN 1842-6298 (electronic), P. 99 – 109, Vol. 4, (2009).
- [16] Ehsanes, A.K.Md. Saleh. "Ridge Regression Estimation Approach to Measurement Error Model". Carleton University, Ottawa (CANADA), Department of Mathematics & Statistics, Indian Institute of Technology, Kanpur - 208 016, (INDIA).
- [17] [https:// STAT 897D](https://STAT897D) - Applied Data Mining and Statistical Learning. "Lesson 5: Regression Shrinkage Methods", The Pennsylvania State University Privacy and Legal Statements Department of Statistics online Programs. Copyright 2015. Last revised 26 / 8 / 2015.
- [18] Shedden, Kerby. "Prediction". Department of Statistics, University of Michigan. November 3, 2014.
- [19] <https://onlinecourses.science.psu.edu/stat512/node/217>. "Applied Linear Models Topic 5a. Ridge Regression (Section 11.2)". Last revised 31 / 8 / 2015.
- [20] Marquardt, Donald W. and Ronald D. Snee Source. "Ridge Regression in Practice". The American Statistician, Vol. 29, Pa. 3-20. No. 1, (Feb., 1975).

