

# EFFECT OF CLUSTERING DATA IN IMPROVING MACHINE LEARNING MODEL ACCURACY

SAMIH M. MOSTAFA<sup>1</sup>, HIROFUMI AMANO<sup>2</sup>

<sup>1</sup>Mathematics Department, Faculty of Science, South Valley University, Qena, Egypt

<sup>2</sup>Research Institute for Information Technology, Kyushu University, Japan

E-mail: <sup>1</sup> samih\_montser@sci.svu.edu.eg, <sup>2</sup> amano@cc.kyushu-u.ac.jp

## ABSTRACT

Supervised machine learning algorithms consider the relationship between dependent and independent variables rather than the relationship between the instances. Machine learning algorithms try to learn the relationship between the input and output from the historical data in order to attain precise predictions about unseen future. Conventional foretelling algorithms are usually based on a model learned and trained from historical data. The instances in the historical data may vary in its characteristics. The variation may be a result of difference in case's pertinence degree to some cases compared to others. However, the problem with such machine learning algorithms is their dealing with the whole data without considering this variation. This paper presents a novel technique to the trained model to improve the prediction accuracy. The proposed method clusters the data using K-means clustering algorithm, and then applies the prediction algorithm to every cluster. The value of K which gives the highest accuracy is selected. The authors performed comparative study of the proposed technique and popular prediction methods namely Linear Regression, Ridge, Lasso, and Elastic. On analysing on five datasets with different sizes and different number of clusters, it was observed that the accuracy of the proposed technique is better from the point of view of Root Mean Square Error (RMSE), and coefficient of determination ( $R^2$ ).

**Keywords:** Prediction accuracy, K-means, clustering, regression, machine learning algorithms.

## 1. INTRODUCTION

Machine learning (ML) is considered as the important subfield of artificial intelligence and is being adopted for numerous of various applications [1, 2]. ML addresses the study and construction of models capable of learning from the data. Understanding *what* and *how* the ML algorithm is learning is an issue for the developers of the ML applications [3]. ML can be classified into unsupervised and supervised. *Unsupervised* learning groups the data into categories depending on the basis of the similarities between data in each group. On the other hand, *supervised* learning means that the machine learns with the assistance of the labeled training data. Estimating unknown (independent, dependent) mapping of a system using a specific number of (independent, dependent) samples is called learning [4, 5]. This process of estimating needs data collected (i.e., training data), and an algorithm that deals with this data and learns from it. Generally speaking, the learning algorithm learns pattern in the data on hand and create a set of rules to map input/output

relation. Data is categorized into labeled (with outcome) or unlabeled (without outcome). Outcome variable(s) may be continuous or distinct, regression is a way of predicting for continuous outcome, and classification is a way of predicting for distinct outcome (i.e., the response to be predicted is the probability or the true of an event/class), the number of classes can be two or more. On the other hand, clustering is applied to unlabeled data using the similarities between observations to group them into clusters [6–8]. The statistical method depends on the characteristics of the data (e.g., similarities between instances in the clustering technique), in other words, the more similarities the better statistical method accuracy. Regression is one of the most common statistical processes for estimating dependent/independent relationship when the dependent to be predicted is a continuous value. The regression line is a refined outline of averages and is drawn in such a way as to reduce the error of the fitted values in relation to the actual values. Equation of the simple linear regression can be defined by the following form:

$$d = C_0 + C_1 I$$

where  $C_0$  is the intercept and the  $C_1$  is the slope of the regression line. In addition to explaining the relationship between dependent  $d$  and independent  $I$  variables, the model also predicts the value of dependent variable from the independent variable values from the equation:

$$\hat{d} = \hat{C}_0 + \hat{C}_1 I$$

the hat symbol  $\hat{\phantom{x}}$  refers to the predicted value of the unknown coefficient/variable. In simple regression, for a given dependent variable there is one independent variable. However, in real cases there is more than one independent variable, so existing of multiple independent variables is called multivariate regression. The mathematical notation can take the form:

$$d = C_1 I_1 + C_2 I_2 + C_3 I_3 + \dots + C_m I_m$$

where  $m$  is the number of independent variables.

### 1.1. Goodness of the Model

The performance of the model can be evaluated using following metrics:

- R-squared: how well the model fits the data.
- Root Mean Square Error (RMSE): how close the estimated values are to the actual values.

### 1.2. Paper Contributions & Novelty

This paper gives a brief discussion of machine learning types, proposes a method for improving the prediction accuracy, and compares between the proposed method and the common methods. In contrast to other techniques, e.g., Clustering Lasso (CL) which selects groups of variables that have the same mechanism of predicting the dependent variable [9], the novelty of the proposed work boils down in benefiting from similarities between instances and applying the selected prediction algorithm for each cluster. The proposed work does not depend on collinearity among variables or the number of variables compared to instances.

### 1.3. Organization

The rest of this paper is organized as follows: Section 2 presents the proposed method

supported with an illustrative example; dataset specification, the prediction algorithms used in the comparison, the parameters used in the comparison, results and discussion are discussed in section 3.

## 2. PROPOSED METHOD

In the proposed method, the authors avail from the similarity attribute of the data by clustering it into groups (clusters) before applying the statistical method. In the proposed method, the number of clusters is determined by the elbow method (heuristic method of validation and interpretation of symmetry within cluster analysis). The data is clustered using K-means algorithm in which the number of clusters resulted from the elbow method is used for the clustering. The selected prediction algorithm is applied for every cluster. For deeper clarification, the next subsection discusses an illustrative example.

Algorithm:	Proposed method
1.	Input data
2.	Select the prediction algorithm
3.	Find the value of K using elbow method
4.	Cluster the data using K-means
5.	Apply the selected prediction algorithm in each cluster

In the proposed method, clustering preprocessing step is applied to the data before applying the prediction algorithm for the purpose of improving the accuracy of model generated by the user-selected prediction algorithm.

### 2.1 Illustrative Example

To clarify the proposed method, an artificial data with 600 observations has been generated, and linear regression algorithm is applied to the data (Figure 1). K-means algorithm is applied to clustering the data, the value of K is determined using elbow method (Figure 2), the data is clustered into three clusters (Figure 3), then linear regression algorithm is applied to the data in each cluster (Figure 4).

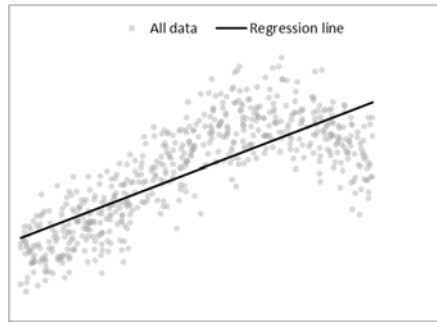


Figure 1. Regression line for all data

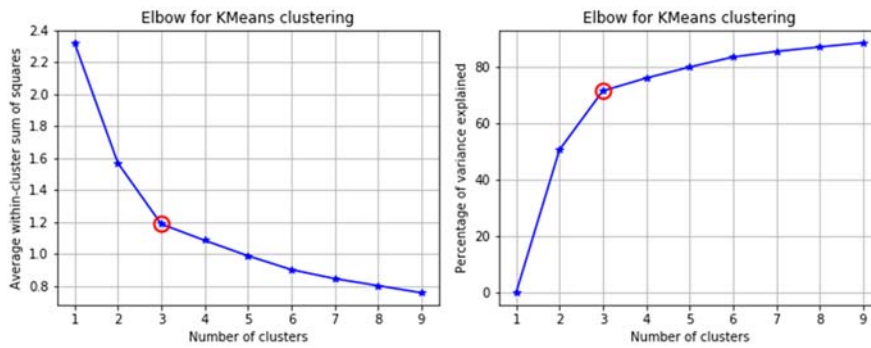


Figure 2. Determining the value of k using elbow method.

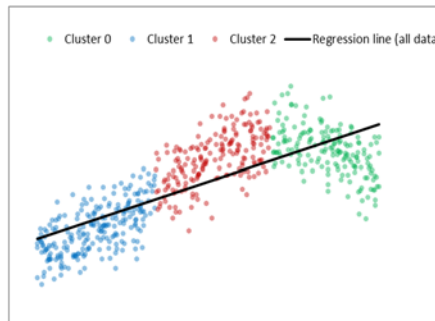


Figure 3. Three clusters of the data

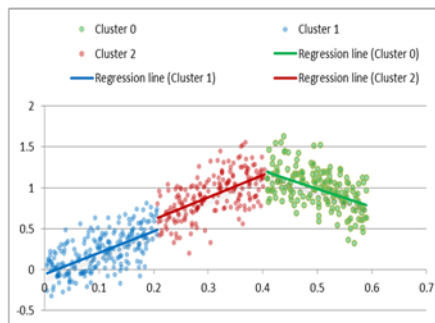


Figure 4. Three regression lines for the clusters

### 3. EXPERIMENTAL IMPLEMENTATION

#### 3.1. Datasets

Six datasets that are commonly used in databases repository are used in the comparative study (Table 1).

#### 3.2. Prediction Algorithms

The comparison is done between four common prediction algorithms namely multiple linear regression (MLR), Ridge, Lasso, and ElasticNet and the proposed method. Table 2 gives short descriptions of the prediction algorithms used.

Table 1. Datasets specifications

Dataset name	#Instances	#Features	References
Diabetes	442	11	[10]
Graduate Admissions	500	8	[11]
California	20640	9	[12]
Diamonds	53940	10	[13]
Boston	506	12	[14]
Iris	150	6	[15]

Table 2. Packages and functions used

Prediction method	Short description	References
MLR	minimizes the residual sum of squares between the observed responses in the dataset, and the responses predicted by the linear approximation.	[16]
Ridge	solves a regression model where the loss function is the linear least squares function, and imposes a penalty on the size of the coefficients.	[17]
Lasso	estimates sparse coefficients. Coefficients that add lightweight value to the model will be zero	[18]
ElasticNet	allows for learning a sparse model where few of the weights are non-zero like Lasso, while still maintaining the regularization properties of Ridge.	[19]

#### 3.3. Performance Measure

The comparisons were done from the point of view of the following parameters:

- **Root mean squared error (RMSE):** indicates how close the forecasted values are to the actual values; hence the lower value of RMSE, the good of the model performance [20]. The mathematical notation can be written as:

$$RMSE = \sqrt{\frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n}}$$

where  $y_i$  and  $\hat{y}_i$  are the actual value and forecasted value of the  $i$ -th observation respectively, and  $n$  is the number cases.

- **Coefficient of determination ( $R^2$  score):** it is a measure of how perfectly the evaluated regression line of the model adapts the data distribution [21]. It can be written as:

$$R^2(y, \hat{y}) = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2};$$

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

#### 3.4. Experimental Results and Discussion

The experiments are conducted on a computer equipped with 16 GB of RAM, Intel core i5-2400 (3.10 GHz), 1 TB of HDD, Gnu/Linux Fedora 28 of OS, and Python (version 3.7) of programming language. Table 3 summarizes the observations by comparing the improvements of the proposed approach versus common four algorithms from the point of view of  $R^2$  score and RMSE. The notable observations are:

- Clustered MLR is better than MLR in  $R^2$  score and RMSE for all datasets.
- Clustered Ridge is better than Ridge in  $R^2$  score and RMSE for all datasets.
- Clustered Lasso is better than Lasso for Iris, Diamond, and Diabetes, and behaves

somewhat similar to Lasso for Admission, Boston, and California.

- Clustered ElasticNet is better than ElasticNet for Iris, Diamond, and Diabetes, and behaves somewhat similar to ElasticNet for Admission, Boston, and California

Table 4 shows the comparisons between the four common methods and the proposed method from the point of view of  $R^2$  score and RMSE.

Table 3. Proposed approach versus MLR, Ridge, Lasso, and ElasticNet

dataset	MLR		Ridge		Lasso		ElasticNet	
	$R^2$ score	RMSE	$R^2$ score	RMSE	$R^2$ score	RMSE	$R^2$ score	RMSE
Admission	better	better	better	better	worst	worst	worst	same
Boston	better	better	better	better	better	better	same	same
California	better	better	better	better	worst	worst	worst	worst
Diabetes	better	better	better	better	better	better	better	better
Diamond	better	better	better	better	better	better	better	better
Iris	better	better	better	better	better	better	better	better

#### 4. Conclusion and Future Works

This paper introduces a method that aims to improve the prediction accuracy of the model by clustering the data and applying the selected algorithm, which is a user choice, for each cluster. Unlike the traditional supervised algorithms which find the relationship between the dependent and independent variables, the proposed approach benefits from the similarities between the instances

to improve the prediction accuracy. Four common algorithms are compared with the proposed method, the results showed that the proposed method achieves significant improvement from the point of view of RMSE, and coefficient of determination  $R^2$ . In the future research avenues, the proposed approach will be analysed in more dataset, other standard error metrics will be considered (e.g., P-value and T-value).

Table 4. Comparisons between common algorithms and the proposed

dataset	Cluster ID	# obs. per cluster	Clustered MLR		MLR		Clustered Ridge		Ridge	
			R <sup>2</sup> score	RMSE	R <sup>2</sup> score	RMSE	R <sup>2</sup> score	RMSE	R <sup>2</sup> score	RMSE
Admission (500 obs.)	0	80	0.90058 7	0.04269 2	0.90308 2	0.04215 3	0.89668 7	0.04352 2	0.90232 9	0.04231 9
	1	85	0.76472 8	0.06693 4	0.7314 8	0.07151 8	0.77273 2	0.06578 5	0.73459 2	0.07109 1
	2	84	0.90039 3	0.03267 7	0.77202 2	0.04943 6	0.90909 8	0.03121 8	0.77161 8	0.04948 8
	3	83	0.87173 9	0.03588 5	0.87628 5	0.03524 6	0.86949 5	0.0362 5	0.87580 2	0.03531 4
	4	82	0.91456 1	0.04097 6	0.91714 2	0.04035 3	0.90643 1	0.04288 2	0.91548 2	0.04075 5
	5	86	0.88446 7	0.06312 7	0.85382 7	0.07100 4	0.88819 3	0.06209 9	0.85386 6	0.07099 4
	Average		<b>0.87274</b> <b>3286</b>	<b>0.04704</b> <b>9204</b>	<b>0.84229</b> <b>3064</b>	<b>0.05161</b> <b>8251</b>	<b>0.87377</b> <b>1396</b>	<b>0.04695</b> <b>0917</b>	<b>0.84228</b> <b>0112</b>	<b>0.05165</b> <b>8971</b>
Improve ment				<b>0.03615</b> <b>1576</b>	<b>0.08851</b> <b>6127</b>			<b>0.03738</b> <b>8136</b>	<b>0.09113</b> <b>7201</b>	
Boston (506 obs.)	0	185	0.81549 9	3.80968 1	0.75999 7	4.34508 4	0.80465 1	3.92008 4	0.75521 4	4.38816 4
	1	137	0.60816 1	4.94037 1	0.32914 5	6.46426 8	0.59972 1	4.99328 5	0.30788 6	6.56589 4
	2	184	0.79888 1	4.06712 9	0.77684 2	4.28417 4	0.80702 5	3.98392 5	0.77468 4	4.30484 3
	Average		<b>0.74084</b> <b>6565</b>	<b>4.27239</b> <b>3928</b>	<b>0.62199</b> <b>4708</b>	<b>5.03117</b> <b>545</b>	<b>0.73713</b> <b>2497</b>	<b>4.29909</b> <b>7932</b>	<b>0.61259</b> <b>4663</b>	<b>5.08630</b> <b>0222</b>
Improve ment				<b>0.19108</b> <b>178</b>	<b>0.15081</b> <b>5953</b>			<b>0.20329</b> <b>5656</b>	<b>0.15476</b> <b>9136</b>	
California (20640 obs.)	0	720	0.66595 8	0.60168 6	0.62772 2	0.63518 9	0.66600 1	0.60164 8	0.62768 9	0.63521 7
	1	113	0.64299 5	0.73097 1	0.60065 6	0.77310 6	0.64299 3	0.73097 3	0.60066 3	0.77309 5
	2	234	0.75353 2	0.44839 8	0.70599 1	0.48973 8	0.72469 8	0.47390 8	0.70594 7	0.48977 5
	3	185	0.64994 4	0.58180 9	0.63040 2	0.59782 9	0.64920 4	0.58242 4	0.63035 5	0.59786 7
	Average		<b>0.67810</b> <b>7166</b>	<b>0.59071</b> <b>6165</b>	<b>0.64119</b> <b>1351</b>	<b>0.62396</b> <b>5652</b>	<b>0.67072</b> <b>195</b>	<b>0.59723</b> <b>8209</b>	<b>0.64116</b> <b>3447</b>	<b>0.62398</b> <b>8367</b>
Improve ment				<b>0.05757</b> <b>3789</b>	<b>0.05328</b> <b>7367</b>			<b>0.04610</b> <b>1354</b>	<b>0.04286</b> <b>9642</b>	
Diabetes (442 obs.)	0	180	0.24455 4	38.0004 9	-0.58888 -	55.1102 7	0.17910 1	39.6125 4	-1.15089 -	64.1204 2
	1	262	0.09290 1	34.3391 5	-1.30065 -	54.6873 9	0.10635 4	34.0836 4	-1.21911 -	53.7096 2
	Average		<b>0.16872</b> <b>7508</b>	<b>36.1698</b> <b>2113</b>	<b>0.94476</b> <b>16</b>	<b>54.8988</b> <b>3005</b>	<b>0.14272</b> <b>5619</b>	<b>36.8480</b> <b>7093</b>	<b>1.18500</b> <b>069</b>	<b>58.9150</b> <b>139</b>
Improve ment				<b>1.17859</b> <b>2683</b>	<b>0.34115</b> <b>4974</b>			<b>1.12044</b> <b>349</b>	<b>0.37455</b> <b>5508</b>	
Diamond (53940 obs.)	0	329	0.87864 9	302.700 3	0.33042 9	711.032 7	0.87838 1	303.034 5	0.32941 9	711.567 8
	1	565	0.40465 4	1986.21 3	-0.51074 -	3164.01 2	0.40428 5	1986.83 8	-0.5104 -	3163.65 3
	2	153	0.59228 7	1007.74 2	0.29129 7	1328.63 7	0.59270 6	1007.22 4	0.29176 9	1328.18 8
	Average		<b>0.62519</b> <b>8361</b>	<b>1098.88</b> <b>5094</b>	<b>0.03699</b> <b>4739</b>	<b>1734.55</b> <b>8121</b>	<b>0.62512</b> <b>3609</b>	<b>1099.03</b> <b>2344</b>	<b>0.03692</b> <b>9956</b>	<b>1734.46</b> <b>9421</b>



ISSN: 1992-8645

[www.jatit.org](http://www.jatit.org)

E-ISSN: 1817-3195

Improve ment		15.8996	0.36647	15.9272	0.36635					
		5597	5484	7726	8189					
Iris (150 obs.)	0	53	0.58589 2	0.15982 4	0.53837 3	0.16874 5	0.60948 6	0.15520 4	0.58034 6	0.16089 1
	1	97	0.50539 3	0.24685 1	0.51907 3	0.24341 4	0.50930 6	0.24587 3	0.51660 3	0.24403 8
	Average		<b>0.54564</b> <b>2513</b>	<b>0.20333</b> <b>7645</b>	<b>0.52872</b> <b>2806</b>	<b>0.20607</b> <b>9331</b>	<b>0.55939</b> <b>6339</b>	<b>0.20053</b> <b>8486</b>	<b>0.54847</b> <b>4215</b>	<b>0.20246</b> <b>4357</b>
	Improve ment				<b>0.03200</b> <b>1092</b>	<b>0.01330</b> <b>4032</b>			<b>0.01991</b> <b>3651</b>	<b>0.00951</b> <b>2148</b>

dataset	Cluster ID	# obs. per cluster	Clustered Lasso		Lasso		Clustered ElasticNet		ElasticNet	
			R <sup>2</sup> score	RMSE	R <sup>2</sup> score	RMSE	R <sup>2</sup> score	RMSE	R <sup>2</sup> score	RMSE
Admission (500 obs.)	0	80	0.39724 8	0.10512 3	0.10319 2	0.12822 7	0.624298 825	0.08299 5	0.476870 26	0.09793 4
	1	85	0.41897 5	0.10518 6	0.21894 6	0.12195 5	0.558337 418	0.09170 8	0.500034 743	0.09757 3
	2	84	- -	0.10425 7	0.24701 5	0.08984 5	0.590989 869	0.06621 7	0.636564 741	0.06241 9
	3	83	- -	0.10233 3	0.28481 6	0.08474 3	0.383689 678	0.07866 7	0.595595 155	0.06372 4
	4	82	0.21780 2	0.12398 3	0.29519 4	0.11769 -	0.653176 12	0.08255 8	0.645685 84	0.08344 5
	5	86	0.25962 2	0.15979 9	0.29872 -	0.15552 3	0.649727 827	0.10991 4	0.667472 07	0.10709 3
	Average		<b>0.20613</b> <b>7182</b>	<b>0.11678</b> <b>0231</b>	<b>0.24131</b> <b>3579</b>	<b>0.11633</b> <b>0496</b>	<b>0.576703</b> <b>29</b>	<b>0.08534</b> <b>2992</b>	<b>0.587037</b> <b>135</b>	<b>0.08536</b> <b>4642</b>
Improve ment				<b>0.14577</b> <b>048</b>	<b>0.00386</b> <b>601</b>			<b>0.017603</b> <b>393</b>	<b>0.00025</b> <b>3615</b>	

Boston (506 obs.)	0	185	0.50421 1	6.24507 9	0.63127 6	5.38567 8	0.555916 205	5.91047 1	0.635021 837	5.35825 1
	1	137	0.49127 7	5.62918 9	0.47095 2	5.74053 7	0.489438 342	5.63935 3	0.473984 885	5.72406 1
	2	184	0.78851 3	4.17063 9	0.64288 7	5.41955 8	0.705213 571	4.92396 2	0.656423 965	5.31584 5
	Average		<b>0.59466</b> <b>7172</b>	<b>5.34830</b> <b>2053</b>	<b>0.58170</b> <b>4998</b>	<b>5.51525</b> <b>7787</b>	<b>0.583522</b> <b>706</b>	<b>5.49126</b> <b>1589</b>	<b>0.588476</b> <b>895</b>	<b>5.46605</b> <b>2459</b>
Improve ment				<b>0.02228</b> <b>3071</b>	<b>0.03027</b> <b>161</b>			<b>0.008418</b> <b>664</b>	<b>0.00461</b> <b>194</b>	

California (20640 obs.)	0	720	0.19451 8	0.93432 8	0.28420 8	0.88077 1	0.398249 064	0.80756 6	0.422555 231	0.79108 8
	1	113	0.34741 2	0.98828 5	0.29673 2	1.02594 3	0.450214 811	0.90710 8	0.436138 572	0.91864 7
	2	234	- -	0.90696 6	0.31536 8	0.74732 9	0.320412 819	0.74457 1	0.486811 721	0.64702 6
	3	185	0.06335 1	0.95170 1	0.28480 6	0.83161 9	0.325355 782	0.80769 9	0.420265 792	0.74873 2
Average		<b>0.14922</b> <b>8941</b>	<b>0.94531</b> <b>9905</b>	<b>0.29527</b> <b>854</b>	<b>0.87141</b> <b>5265</b>	<b>0.373558</b> <b>119</b>	<b>0.81673</b> <b>5984</b>	<b>0.441442</b> <b>829</b>	<b>0.77637</b> <b>343</b>	
Improve ment				<b>0.49461</b> <b>637</b>	<b>0.08480</b> <b>99</b>			<b>0.153779</b> <b>165</b>	<b>0.05198</b> <b>858</b>	



	0		0.01644	43.3596	-	69.1600	0.000281	43.7145	3.415714	91.8730
		180	9	7	1.50228	3	698	8	253	2
Diabetes (442 obs.)	1		-	36.4817	-	58.4636	0.023699	36.4794	2.644023	68.8260
		262	0.02383	1	1.62935	8	919	7	327	4
	Average		<b>0.00368</b>	<b>39.9206</b>	<b>1.56581</b>	<b>63.8118</b>	<b>0.011709</b>	<b>40.0970</b>	<b>3.029868</b>	<b>80.3495</b>
			<b>829</b>	<b>8967</b>	<b>185</b>	<b>5592</b>	<b>111</b>	<b>2716</b>	<b>79</b>	<b>3014</b>
					<b>0.99764</b>	<b>0.37440</b>			<b>0.996135</b>	<b>0.50096</b>
					<b>4487</b>	<b>0116</b>			<b>44</b>	<b>7497</b>
	0		0.87573	306.308	0.32562	713.579	0.743478	440.100	0.320155	-
		329	9	3	2	5	863	9	658	998.397
Diamond (53940 obs.)	1		0.40457	1986.35	-		0.253176	2224.60	2.206191	-
		565	5	3	0.50916	3162.36	236	1	315	4609.33
	2	153	0.59348	1006.25	0.29303	1327.00	0.402987	1219.44	0.416165	1205.91
		23	8	7	3	2	813	9	349	6
	Average		<b>0.62460</b>	<b>1099.63</b>	<b>0.03649</b>	<b>1734.31</b>	<b>0.466547</b>	<b>1294.71</b>	<b>0.703393</b>	<b>2271.21</b>
			<b>0687</b>	<b>9291</b>	<b>6857</b>	<b>3961</b>	<b>637</b>	<b>6959</b>	<b>874</b>	<b>4188</b>
					<b>16.1138</b>	<b>0.36595</b>			<b>1.663280</b>	<b>0.42994</b>
					<b>214</b>	<b>1428</b>			<b>779</b>	<b>502</b>
	0		-	0.25231	-	0.69119	0.032111	0.25231	1.973883	0.42829
		53	0.03211	8	6.74515	4	237	8	919	9
Iris (150 obs.)	1		-	0.38414	-	0.38273	0.197813	0.38414	0.367240	0.27920
		97	0.19781	9	0.18902	7	953	9	443	6
	Average		<b>0.11496</b>	<b>0.31823</b>	<b>3.46708</b>	<b>0.53696</b>	<b>0.114962</b>	<b>0.31823</b>	<b>0.803321</b>	<b>0.35375</b>
			<b>26</b>	<b>3686</b>	<b>467</b>	<b>541</b>	<b>595</b>	<b>3686</b>	<b>738</b>	<b>2519</b>
					<b>0.96684</b>	<b>0.40734</b>			<b>0.856890</b>	<b>0.10040</b>
					<b>1711</b>	<b>7886</b>			<b>97</b>	<b>588</b>

REFERENCES:

- [1] T. Segaran, Programming Collective Intelligence: Building Smart Web 2.0 Applications. 2007.
- [2] D. Vallejo-Huanga, P. Morillo, and C. Ferri, "Semi-Supervised Clustering Algorithms for Grouping Scientific Articles," Procedia Comput. Sci., vol. 108, pp. 325–34, 2017.
- [3] V. K. Ayyadevara, Pro Machine Learning Algorithms. 2018.
- [4] V. Cherkassky and F. M. Mulier, Learning from Data : Concepts, Theory, and Methods / V. Cherkassky, F. Mulier. 1998.
- [5] W. A. Yousef and S. Kundu, "Learning algorithms may perform worse with increasing training set size: Algorithm-data incompatibility," Comput. Stat. Data Anal., vol. 74, pp. 181–97, 2014.
- [6] K. Chen and L. Liu, "The 'Best K' for Entropy-based Categorical Data Clustering.," 2005, pp. 253–62.
- [7] M. N. Murty, P. J. Flynn, and A. K. Jain, "Data Clustering: A Review," ACM Comput. Surv., vol. 31, p. 43, 1999.
- [8] A. K. Jain and R. C. Dubes, "Algorithms for Clustering data." 1988.
- [9] Q. Yu and B. Li, "Regularization and Estimation in Regression with Cluster Variables," Open J. Stat., vol. 04, no. 10, pp. 814–25, 2014.
- [10] N. D. of Statistics, "Diabetes Data-1-5-2019," NCSU Department of Statistics. [Online]. Available: <https://www4.stat.ncsu.edu/~boos/var.select/diabetes.html>. [Accessed: 01-May-2019].
- [11] M. S. Acharya, "Graduate Admissions-1-5-2019," Kaggle Inc, 2018. [Online]. Available: <https://www.kaggle.com/mohansacharya/graduate-admissions>. [Accessed: 01-May-2019].
- [12] C. Nugent, "California Housing Prices-3-5-2019," Kaggle Inc, 2017. [Online]. Available: <https://www.kaggle.com/camnugent/california-housing-prices>. [Accessed: 03-May-2019].
- [13] M. Shiva, "Diamonds-20-4-2019," Kaggle Inc, 2017. [Online]. Available:





- <https://www.kaggle.com/shivam2503/d/amonds>. [Accessed: 20-Apr-2019].
- [14] I. Kaggle, “Boston Housing-25-4-2019.” [Online]. Available: <https://www.kaggle.com/c/boston-housing>. [Accessed: 25-Apr-2019].
- [15] I. Kaggle, “Iris Dataset-25-4-2019.” [Online]. Available: <https://www.kaggle.com/jchen2186/machine-learning-with-iris-dataset>. [Accessed: 25-Apr-2019].
- [16] D. W. Gareth James Trevor Hastie, Robert Tibshirani, An introduction to statistical learning : with applications in R. New York : Springer, [2013] ©2013.
- [17] Scikit-learn, “Ridge Regression-26-4-2019.” [Online]. Available: [https://scikit-learn.org/stable/modules/generated/sklearn.linear\\_model.Ridge.html](https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.Ridge.html). [Accessed: 26-Apr-2019].
- [18] Scikit-learn, “Lasso Regression-27-4-2019.” [Online]. Available: [https://scikit-learn.org/stable/modules/linear\\_model.html](https://scikit-learn.org/stable/modules/linear_model.html). [Accessed: 27-Apr-2019].
- [19] Scikit-learn, “Elastic Net Regression-28-4-2019.” [Online]. Available: [https://scikit-learn.org/stable/modules/linear\\_model.html#elastic-net](https://scikit-learn.org/stable/modules/linear_model.html#elastic-net). [Accessed: 28-Apr-2019].
- [20] I. B. Aydilek and A. Arslan, “A hybrid method for imputation of missing values using optimized fuzzy c-means with support vector regression and a genetic algorithm,” *Inf. Sci. (Ny)*, vol. 233, pp. 25–35, 2013.
- [21] T. Beysolow II, *Introduction to Deep Learning Using R*. Apress, 2017.