1 **COMPREHENSIVE CLUSTERWISE LINEAR REGRESSION FOR PAVEMENT**
2 **MANAGEMENT SYSTEMS**
3

4

5 **Mukesh Khadka**
6 Ph.D. Candidate
7 Cell: (702) 683-2722
8 Email: khadkam@unlv.nevada.edu
9
10 **Alexander Paz, Ph.D., P.E., Corresponding Author**
11 Associate Professor
12 Director, Transportation Research Center
13 Office: (702) 895-0571
14 Cell: (702) 688-3878
15 Fax: (702) 895-3936
16 Email: apaz@unlv.edu
17 http://web.unlv.edu/centers/trc/paz/
18
19 Civil and Environmental Engineering and Construction
20 University of Nevada, Las Vegas
21 4505 Maryland Parkway, PO Box 454007, Las Vegas, NV 89154-4007
22
23
24 **ABSTRACT**

25 A comprehensive mathematical program was formulated to determine simultaneously 1) an
26 optimum number of pavement clusters, 2) cluster memberships of pavement samples, 3) cluster-
27 specific significant explanatory variables, and 4) estimated regression coefficients for Pavement
28 Performance Models (PPMs). Simulated Annealing coupled with All-Subset Regression was
29 proposed to solve the mathematical programming. The proposed algorithm was capable to
30 identify and address potential multicollinearity issues. All possible combinations of the
31 explanatory variables were examined to select the best model that provided a balance among 1)
32 the number of PPMs; 2) the number of explanatory variables; 3) the resources required to
33 develop, maintain, and use these models; and 4) the explanatory power. For the dataset used in
34 this research, 6-cluster models were determined as part of the optimum solution. The predictive
35 capabilities of the resultant models were investigated, and results showed that the models
36 provided few prediction errors without any overfitting issues.

## INTRODUCTION

Pavement deteriorates over time due to the combined effects of traffic and environmental factors. To keep pavement in a serviceable condition, highway agencies primarily have two alternatives: 1) permit the pavement to deteriorate until its condition falls below the serviceability limit, and then perform rehabilitation or reconstruction work; or 2) intervene with the deterioration by performing a series of maintenance activities that retard the deterioration process and essentially delay the type of substantial failure requiring major rehabilitation or reconstruction.

Considering that a typical cost of the maintenance is 15% to 20% of the cost for rehabilitation or reconstruction (Hajj et al. 2010), agencies are more focused on preserving and maintaining existing facilities (Davies and Sorenson 2000; Labi and Sinha 2003). However, the challenge is to find the pavement segments that require maintenance as well as appropriate times to execute such activities. Hence, there is a need to develop a proactive approach to identify potential pavement segments for improvement. Pavement performance models (PPMs) – one of several critical components required to achieve this proactive approach – seek to capture historical patterns of pavement deterioration that can be used to estimate an appropriate time for maintenance so that the condition of a pavement can be improved before a serviceability limit is reached.

In practice, it is very important to achieve a balance among the number of PPMs; the number of explanatory variables; the resources required to develop, maintain, and use these models; and the associated explanatory power. To determine this balance, PPMs typically are developed by using clusters of pavement samples. Instead of estimating the cluster memberships by using statistical methods, a few predefined explanatory variables are used to assign pavement

59    samples into clusters. In terms of performance, clusters formed in this way likely include

60    heterogeneous pavement samples.

61         The existing state-of-the-art methods propose Clusterwise Linear Regression (CLR) to

62    determine pavement clusters and associated PPMs simultaneously, using a single objective

63    function. In CLR, various clusters are formed so that homogenous pavement samples, in terms of

64    the effects of the explanatory variables on the dependent variable of a present regression model,

65    are assigned within a cluster (Park et al. 2015). The homogeneity of pavement samples in a

66    cluster is defined by the effects of the observed values of explanatory variables on the estimated

67    dependent variable, the Present Serviceability Index (PSI), by the regression model.

68    Observations of all the pavement samples assigned to a cluster fit the same PPM such that the

69    overall sum of squared errors (SSE) within clusters is minimal.

70         CLR first was implemented by Spath (1979) for data partition and estimation of

71    regression models within each cluster, simultaneously. The approach has been expanded further,

72    and implemented in many studies (DeSarbo et al. 1989; Wedel and SteenKamp 1989; Lau et al.

73    1999; Carbonnea et al. 2011; Schlittgen 2011; Zhen et al. 2012; Tan et al. 2013; Lu et al. 2014).

74    However, in the field of pavement management, to the best knowledge of the authors, only four

75    studies (Luo and Chou 2006; Luo and Yin 2008; Zhang and Durango-Cohen 2014) have been

76    performed using CLR.

77         In a recent study (Zhang and Durango-Cohen 2014), CLR with multiple explanatory

78    variables was proposed to account for heterogeneity in pavement deterioration. The study used

79    the data collected during the AASHO Road Test (Highway Research Board 1962), which is no

80    longer the best available data nor representative of existing conditions. This data was collected at

81    a single site, and over 50 years ago, when materials and construction techniques were different.

82   The study estimated models with the objective of minimization of the residual sum of squares

83   (RSS). The number of models were determined subjectively using the trends of RSS and Akaike

84   Information Criteria (AIC) over the number of clusters. In addition, the study investigated the

85   presence of overfitting in the CLR models, using a procedure proposed by Brusco et al. (2008).

86   In this current study, overfitting means that most of the variations in the dependent variable

87   appears to be explained by the estimated model; however, the actual relationship between the

88   dependent variable and some of the explanatory variables and/or the functional form of the

89   model is not really captured. Overfitting typically is evidenced during validation when the model

90   is used to estimate values for the dependent variable, using data that was not used for model

91   development. Later in this paper, the section on Model Performance provides a rigorous

92   explanation of a procedure to determine potential overfitting in a model.

93        To address some of the limitations of previous models, a mathematical programming

94   framework within the CLR approach is proposed to determine simultaneously the optimal

95   number of clusters, the assignment of segments into clusters, and the associated PPMs (Khadka

96   and Paz, 2017b). In this study, the Bayesian Information Criteria (BIC) (Schwarz 1978) was used

97   as the objective function. BIC penalizes more for the inclusion of additional parameters than

98   does AIC (Kadane and Lazar 2004). On the other hand, several studies showed that the number

99   of parameters in a model selected using AIC was overestimated (Geweke and Meese, 1981; Katz,

100   1981; Koehler and Murphree, 1988; Kadane and Lazar 2004).

101        BIC is one of the most popular log-likelihood-based information criteria used for model

102   selection. As BIC is an increasing function of the error variance and free parameters to be

103   estimated, minimizing BIC reduces unexplained variations in the dependent variable, the number

104   of explanatory variables, or both (Uzoma and Jeremiah, 2016). In case of a large sample size,

105    BIC is consistent in the sense that the probability of the selected model being the true model

106    approaches '1' (Rao and Wu 1989; Yang 2005; Maydeu-Olivares and García-Forero 2010; Vrieze

107    2012, Kim et al. 2012).

108         In addition, the proposed framework tests the significance of explanatory variables. To

109    the best of the authors' knowledge, all the existing literature about pavement management and

110    PPMs estimation using CLR suffers from the limitation that variables included in the PPMs are

111    assumed to be significant. However, the effects of variables without any evidence of significance

112    can affect clustering and regression analyses. Therefore, heterogenous samples can be assigned

113    together erroneously (Fowlkes et al. 1988); therefore, it becomes challenging to discover the

114    underlying pavement clusters that exhibit similar performance behavior (Gupta and Ibrahim

115    2007).

116         This problem is illustrated in Figure 1, using data from the Pavement Management

117    System (PMS) of the Nevada Department of Transportation (NDOT). In this example, 54

118    randomly selected pavement samples were considered. Each pavement sample was represented

119    by a dependent variable, PSI, and two explanatory variables, Age and Average Daily Traffic

120    (ADT).

121         The variables PSI and Age had a significant linear relationship (p-value = 0.001), as

122    shown in Figure 1a. The estimated BIC and root mean square error (RMSE) for the model were

123    85 and 0.2916, respectively. However, the relationship between PSI and ADT was not clear, as

124    shown in Figure 1b. The estimated BIC and RMSE for the model were 251 and 0.4572,

125    respectively. When both Age and ADT were included in the model as explanatory variables, the

126    estimated BIC was increased to 90, with a slight decrease in RMSE by 0.0003. Hence, if an

127    irrelevant variable, ADT in this example, is included in a CLR analysis without checking its

128   significance, it increases the BIC. In addition, it causes a loss of efficiency in the model. The

129   estimated clustering and regression models may not capture the correct underlying relationships

130   among the variables when a variable is included in the model without sufficient evidence of its

131   significance.

132        Assignment of pavement samples into clusters using predefined and fixed explanatory

133   variables, instead of estimation, introduces bias into the statistical analysis (Gupta and Ibrahim

134   2007). The available data are not fully utilized for clustering because the performance behavior

135   represented by historical PSI is ignored. In addition, clustering using explanatory variables that

136   do not provide any information about the underlying clustering structure does not reveal the

137   underlying cluster assignments.

138        A legitimate assignment of pavement samples into homogeneous clusters to minimize the

139   estimation error can be obtained using the relevant explanatory variables that exhibit the

140   strongest effects on the dependent variable (Fowlkes et al. 1988; Liu and Ong 2008; and Maugis

141   et al. 2009). The strength of the effects of explanatory variables on the dependent variable often

142   is assessed by comparing p-values with the desired level of significance ($\alpha$). A p-value represents

143   the significance of the estimated coefficient for an explanatory variable. If the p-value for an

144   explanatory variable is greater than $\alpha$, there is not enough evidence to claim that the estimated

145   coefficient is likely to be different from zero. In other words, changes in the explanatory variable

146   do not reflect changes in the dependent variable. Hence, such explanatory variables having p-

147   values greater than the desired $\alpha$ usually are excluded from the model during model estimation

148   process.

149        A variable selection procedure can be utilized to select the best subset of potential

150   explanatory variables. This procedure must distinguish between relevant and irrelevant variables

151 in order to provide the best regression models. Typically, the fewest number of explanatory

152 variables that sufficiently explain most of the variances in the dependent variable are selected as

153 the best model specification. In terms of data analysis and statistics, numerous methodologies for

154 variable selection are available in the literature (Thompson 1978; Tibshirani 1996; Baumann

155 2003; Efron et al. 2004; Mehmood et al. 2012; Brusco 2014). In this study, the All-Subset

156 Regression procedure (Garside 1965; Gorman and Toman 1966; Hocking and Leslie 1967;

157 Mallows 1973; Berk 1978; Efron et al. 2004) was used to select variables for CLR analysis. All

158 ($2^P$-1) possible subsets of potential explanatory variables, $P$, were examined. BIC was used as a

159 criterion for comparing models with different subsets of variables.

160 It is not recommended to use least squares estimation and variable selection techniques

161 under the presence of multicollinearity (Gunst and Webster 1975). Strongly-correlated clustering

162 variables may overweight one or more underlying constructs and produce loss in efficiency

163 (Ketchen and Shook 1996). Typically, multicollinearity inflates the variance of regression

164 parameters and makes correct identification of significant variables challenging (Abdul-Wahab et

165 al. 2005; Dorman et al. 2013; Ohlemüller et al. 2008). However, strongly correlated variables

166 may not be a problem in all cases (Harrell 2001). In addition, if the collinearity between two

167 variables remains constant, their estimated parameters are likely to have low standard errors; the

168 problem would be serious if the standard errors of the correlated variables are high (Washington

169 et al. 2011). The best way to address multicollinearity is to conduct a carefully designed

170 experiment that considers the trade-off between removing and keeping potential explanatory

171 variables that are expected to cause multicollinearity. Judgement and iterations are required to

172 determine the best model specification that minimizes the effects of multicollinearity

173 (Washington et al. 2011).

174     This study investigated the effects of highly-correlated explanatory variables. The

175     Variance Inflation Factor (VIF), used to examine potential issues due to multicollinearity

176     (Marquardt 1970; Mansfield and Helms 1982), is defined as $1/(1 - R_i^2)$, where $R_i^2$ is the R$^2$ for

177     an explanatory variable, $X_i$ regressed on the remaining explanatory variables. When no

178     explanatory variables are correlated, the VIF is equal to '1'. As the degree of collinearity

179     increases, both the variance of regression coefficient and the VIF increase (Yoo et al. 2014). Tacq

180     (1997) showed that large VIF is an indicator of multicollinearity. In general, a VIF greater than

181     '10' is considered unacceptable (Neter et al. 1996; Midi et al. 2010), even though no formal rule

182     exists in the literature.

183     To avoid prespecifying the significance of potential explanatory variables, this paper

184     proposes a comprehensive CLR framework that determines, simultaneously, the optimal number

185     of pavement clusters, the assignment of segments into clusters, and the corresponding PPMs

186     using only likely significant explanatory variables. The proposed framework simultaneously

187     seeks for 1) the optimal number of clusters, 2) the combination of significant explanatory

188     variables that provides the best goodness of fit, and 3) assigns segments into clusters. In the

189     study, the likely significance of the explanatory variables was tested for each cluster model;

190     hence, different clusters may include different significant explanatory variables.

191     Considering the simultaneous and extensive search for significant explanatory variables

192     and the optimal number of clusters, the PPMs developed under the proposed framework were

193     expected to provide superior explanatory power compared to existing approaches. The proposed

194     framework was tested using pavement data from the entire State of Nevada. The results illustrate

195     the advantage of solving simultaneously for the three types of parameters listed above.

196 **METHODOLOGY**

197 **Problem formulation**

198 This section describes a mathematical program that was formulated to describe the proposed

199 CLR problem. Among various pavement performance measures available in the literature, PSI is

200 a widely accepted measure that serves as a unified standard to measure pavement serviceability

201 (Shoukry et al. 1997; Terzi 2006; Attoh-Okine and Adarkwa 2013). PSI is understood easily by

202 both road users and legislators (Hudson et al. 2015). This study used PSI as the dependent

203 variable, *y*. Multiple linear regression PPMs were estimated with functional form expressed by:

204
$$y_{it} = \beta_{0k} + \sum_{j=1}^{J} \beta_{jk} * x_{ijt}$$
(1)

205 The objective function was to minimize BIC, expressed as:

206
$$Min.\ BIC = O + O*ln(2\pi) + O*ln\left(\frac{SSE}{O}\right) + \left(\delta + K\text{-}1\right)*ln(O)$$
(2)

207 where *SSE* is total sum of squared errors, expressed by:

208
$$SSE = \sum_{k=1}^{K} \sum_{i=1}^{I} \sum_{t=1}^{T_i} \left(\beta_{0k} + \sum_{j=1}^{J} \beta_{0k} * x_{ijt} - y_{it}\right)^2 * p_{ik} \ \forall\ i \in I, j \in J, t \in T_i, k \in K$$
(3)

209 and the quantity $(\delta + K\text{-}1)$ is the total number of free parameters to be estimated for *K* clusterwise

210 regression models (DeSarbo and Corn 1988). Intercepts $(\beta_{0k})$, coefficients for cluster-specific

211 significant explanatory variables $(\beta_{jk})$, the optimum number of clusters $(K)$, and cluster

212 memberships $(p_{ik})$ were the decision variables to be determined. In addition, the proposed

213 mathematical programming included the following constraints:

214 Constraints for significant variables:

215
$$\delta = \sum_{k} \sum_{j} v_{jk} \ \forall\ j = 0,...,J,\ k \in K$$
(4)

216
$$v_{jk} = \begin{cases} 1, & if\ \beta_{jk}\ is\ significant; \\ 0, & Otherwise \end{cases} \ \forall\ j = 0,...,J, k \in K$$
(5)

217        Membership constraints:

218        $\sum_k p_{ik}=1 \ \forall i \in I, k \in K$         (6)

219        $p_{ik} = \begin{cases} 1, & \text{if sample } i \text{ is assigned to cluster } k; \\ 0, & \text{Otherwise} \end{cases} \ \forall \ i \in I, k \in K$         (7)

220        Constraints for feasible partitions:

221        $C_k = \{i \mid p_{ik}=1 \forall i \in I, k \in K\}$         (8)

222        $C_{k'} \cap C_{k''} = \text{null} \ \ \forall k' \neq k'', k' \text{ and } k'' \in K$         (9)

223        $\bigcup_{k \in K} |C_k| = I$         (10)

224        $\sum_{i \in C_k} T_i \geq n \ \forall \ C_k$         (11)

225        Constraints for range of clusters:

226        $1 \leq k \leq K_{max}$         (12)

227        $K_{max} = F(I, T_i, n)$         (13)

228

229        The constraint expressed by (4) provided the total number of significant explanatory

230  variables, including intercepts for all the clusters. The sum of elements in each column of the

231  binary matrix, **V**, of size ($J+1$ x $K$) provided the number of significant explanatory variables and

232  an associated intercept for a particular cluster. According to the constraint expressed by (5), the

233  element $v_{jk}$ was equal to '1' if an estimated coefficient ($\beta_{jk}$) was significant in cluster $k$;

234  otherwise, $v_{jk}$ was '0' (Eq. 5). The significance of an explanatory variable as well as an intercept

235  was determined by using the p-value of its estimated regression coefficient.

236        Constraints expressed by (6) and (7) ensured that a pavement sample was assigned

237  exclusively to a single cluster. A binary indicator variable, $p_{ik}$, was used to define the

238    membership of a sample. Indicator $p_{ik}$ equaled '1' if and only if a pavement sample $i$ belonged

239    to cluster $k$. Otherwise, $p_{ik}$ was '0'.

240          The feasibility of the resulting clustering was guaranteed by constraints expressed by (8)

241    - (11). Constraints expressed by (8) – (10) prevented the overlap of members among clusters;

242    that is, pavement samples were divided exclusively into $K$ clusters. Constraint (11) warranted

243    that the number of observations for each cluster was no less than the minimum number of

244    observations, $n$, in order to obtain the statistically reliable estimation of coefficients.

245          Constraints expressed by (12) and (13) were used to prevent a search beyond a feasible

246    number of clusters. If the pavement sample had more than $n$ observations, the sample alone

247    could form a cluster. In reality, none of the pavement samples had more than $n$ observations.

248    Hence, samples were grouped into clusters to provide enough observations. All observations of a

249    sample needed to be assigned to the same cluster.

250          The constraint expressed by (13) denoted the maximum number of feasible clusters. A

251    procedure to calculate this maximum number was denoted by function F (Khadka et al., 2017).

252    The procedure iteratively searched for the best combinations of the pavement samples to form a

253    cluster such that each cluster had the required minimum number of observations. In the first step,

254    it searched pavement samples with $n$ or more observations. In this case, each pavement sample

255    could form a cluster and was assigned to an individual cluster. Once all such cases were

256    searched, the procedure searched two or more pavement samples, where a total number of

257    observations equaled to $n$. In this step, all possible combinations of pavement samples with a

258    total number of observations equal to $n$ were searched to find the maximum number of

259    combinations. No sample could be assigned to more than one cluster. After assigning all possible

260    combinations, the algorithm seeks for the combination(s) of samples having the minimal number

261    of extra observations in each cluster.

262    **Solution to the mathematical program**

263    This study integrated Simulated Annealing (SA) (Krickpatrick et al. 1983; Cěrny, 1985) with

264    Ordinary Least Square (OLS) to solve the proposed mathematical program, which is described as

265    follows by means of algorithmic steps and a discussion regarding the details. SA was chosen

266    because it provides a probabilistic mechanism to seek a global optimum in a large search space

267    that involves discrete variables, such as cluster membership. Thus, SA was used to determine the

268    cluster memberships ($p_{ik}$) of the pavement samples. For each accepted cluster, the VIF for all

269    explanatory variables were calculated as discussed in the introduction. Highly correlated

270    explanatory variables that had VIFs greater than a predefined limiting VIF were excluded. All-

271    subset regressions were utilized to find the best model and to estimate the associated regression

272    coefficients ($\beta_{jk}$). BIC and the level of significance, $\alpha$, were used as the criteria to select the best

273    model. Hence, selected models included only significant explanatory variables at a given $\alpha$.

274          The algorithm utilized to solve the proposed mathematical program is described as

275    follows, and is illustrated in Figure 2.

276    Step 1. Set $K = 2$, $BIC_{min} =$ infinity, and $N = 1$.

277    Step 2. Calculate the maximum number of feasible clusters, $K_{max}$, utilizing function F,

278          described above, as part of the constraint expressed by (13).

279    Step 3. For a given $K$, randomly assign pavement samples into clusters using the following

280          steps:

281       Step 3.1. Generate a random number $u \sim U(1, K)$ and assign it to each of the pavement

282              sample used for the estimation of CLR models. When a sample is assigned to a

283              cluster, all observations associated with that sample are assigned to this cluster.

284       Step 3.2. Find the total number of observations assigned to each of the clusters, (i.e., 1 to

285              $K$).

286       Step 3.3. If all the clusters have at least $n$ observations, then go to Step 4; otherwise, repeat

287              Steps 3.1 and 3.2 until all the clusters have at least $n$ observations. Let $C_K^N \ \forall \ 1 \leq$

288              $k \leq K$ be the valid initial clusters.

289    Step 4. All-subset regression: Repeat the following steps for all $K$ clusters.

290       Step 4.1. Calculate $VIF$ for all explanatory variables. Exclude variables that have $VIF >$

291              $VIF_{max}$. Let $\hat{J}$ be the set of explanatory variables with $VIF < VIF_{max}$.

292       Step 4.2. Generate all possible $2^{|\hat{J}|} - 1$ subsets of $\hat{J}$.

293       Step 4.3. Estimate $\beta_{jk}$ for all subsets, using OLS, and calculate $BIC$ for all the models.

294       Step 4.4. Rank models in ascending order, using $BIC$.

295       Step 4.5. Select the model that has the minimum $BIC$ and all significant explanatory

296              variables with $p\text{-}value < \alpha$.

297    Step 5. Calculate the total number of free parameters to be estimated, $(\delta + K\text{-}1)$. Calculate $BIC$

298       using Eq. 2.

299    Step 6. Using the following steps, generate valid neighborhood clusters near to the previous

300       ones.

301    Step 6.1. Select $N_{ps}$ pavement samples randomly. For each of the selected samples, assign a

302              new membership by generating a random number $u_1 \sim U(1, K)$. If the new

303              membership is the same as previously, regenerate a random number $u_2 \sim U(1, K)$

304                  until a different outcome is obtained. Repeat this process until the memberships of

305                  all selected samples are different from those previously assigned.

306        Step 6.2. If all clusters have at least $n$ observations, go to Step 7; otherwise, repeat Step 6.1.

307                  until all clusters have at least $n$ observations. Let $C_K^{N+1}$ be the new set of valid

308                  neighborhood clusters.

309        Step 7. For $C_K^{N+1}$, repeat Step 4 to estimate $\beta_{jk}$ for all $K$ clusters.

310        Step 8. Calculate the total number of free parameters to be estimated, $(\delta+K\text{-}1)$, and evaluate

311             $BIC_K^{N+1}$, using the Eq. 2.

312        Step 9. Search of a solution.

313          Step 9.1. Calculate $\Delta BIC = BIC_K^{N+1} - BIC_K^{N}$.

314          Step 9.2. Check the following two conditions:

315              a. If $\Delta BIC < 0$ , accept current set of clusters, $C_K^{N+1}$, and the corresponding $\beta_{jk}$; go

316                  to Step 10, otherwise, go to Step b.

317              b. Generate a random number $u''{\sim}U(0,1).$ Calculate the acceptance probability,

318                  $p_{accept} = exp\left(\frac{-\Delta BIC}{B*T}\right),$ where $B$ is the Boltzmann's constant. If $p_{accept} > u'',$

319                  accept the current set of clusters, $C_K^{N+1}$, and the corresponding $\beta_{jk}$. Go to Step 10;

320                  otherwise, return to Step 6.

321        Step 10. Counter and temperature update.

322          Step 10.1. Repeat Steps 6 to 9 for $N_{max}$ times.

323          Step 10.2. If $\theta < \theta_{min}$, stop the algorithm. Otherwise, reduce the temperature by

324                multiplying the current temperature by $\lambda$, set $N =1$, and go to Step 6.

325        Step 11. Stopping criteria.

326          Step 11.1. Update $BIC_{min}$ with the smallest between the values obtained in Step 10 and the

327        current $BIC_{min}$. Set $K_{optimal} = K$.

328        Step 11.2. Repeat Steps 3 to 10 for $K_{max} - 1$ times.

329        To seek a global solution, this algorithm used a probabilistic approach during the search

330    process. The initial solution was improved repetitively by making small changes until a better

331    solution was obtained (Sridhar and Rajendran 1993; Johnson et al. 1989). The algorithm

332    accepted better solutions as well as non-improving (worse) solutions at a certain probability

333    (Dolan et al. 1989; Rutenbar 1989; Aarts et al. 2005). This probability decreased continuously

334    over iterations, and depended on 1) the difference between the BICs of the current solution and a

335    newly selected solution, and 2) the current temperature (Nikolaev and Jacobson 2010).

336        Initially, at a high temperature, the algorithm accepted worse solutions, which caused

337    larger increments in BIC. As the temperature went down, the algorithm accepted worse solutions

338    with relatively smaller increments in BIC. Finally, when the temperature dropped to zero, the

339    algorithm no longer accepted worse solutions. This enabled occasional 'uphill' moves that helped

340    the algorithm to escape from the local minima. Thus, the algorithm tried to explore the entire

341    solution space to seek for a global solution (Dolan et al. 1989). Previous studies have shown that

342    the algorithm converged to a global minimum when an infinitely slow cooling schedule was

343    utilized (Román-Román et al. 2012).

344    **Application of CLR Models**

345    Luo and Chao (2008) proposed a procedure to apply CLR models to estimate pavement

346    conditions. However, the proposed procedure applies only for cases when pavement age is the

347    only independent variable. In addition, the procedure cannot be used to estimate the condition of

348    a pavement sample that was not used to develop the CLR model. In other words, the procedure

349    cannot be used to determine the cluster memberships of the pavement samples that are not

350    included in the estimation process.

351          To address this issues, this study proposed a heuristic to closely assign the cluster

352    membership to a pavement sample. It was assumed that the new sample had observations for all

353    the explanatory variables included in the estimated CLR models as well as the dependent

354    variable (i.e. PSI), for at least one year. The following procedure could be used to estimate PSI

355    using CLR models and the observations for a pavement sample:

356      1. Estimate $\widehat{PSI}_t^k$ separately for all $T$ observations of a sample, using each of the $K$

357          estimated CLR models.

358      2. Calculate the overall sum of squared error (SSE) for each of the models, $\sum_k \Delta PSI_t^k =$

359          $\sum_k \left(\widehat{PSI}_t^k - PSI_t^k\right)^2$

360      3. The sample is assigned to the model associated with the least overall SSE.

361    **EXPERIMENT AND RESULTS**

362    **Data**

363    Experiments were performed using the Pavement Management System (PMS) of the Nevada

364    Department of Transportation (NDOT). The data included condition monitoring and roadway

365    inventory data collected throughout the entire State of Nevada. Potential explanatory variables

366    used in this study could be divided as follows:

367      1. Continuous explanatory variables:

368          • *age* - pavement age since the last M&R treatment;

369          • *adt* - average daily traffic in one direction;

370          • *trucks* - average daily trucks in one direction;

- *elevation* - midpoint elevation of a segment;

- *precip* - average annual precipitation (cm/yr);

- *min_temp* - minimum average annual temperature ($^0$C);

- *max_temp* - maximum average annual temperature ($^0$C);

- *wet_days* - total number of wet days in a year;

- *freeze_thaw* - total number of freeze-thaw cycles that a pavement experienced in a year;

- *rut_depth* - average ride rut depth (cm);

2. Categorical explanatory variables:

- Two dummy variables for the number of lanes were encoded as:

  o *lane ≤ 2* was equal to '1' if the pavement sample had two or less lanes and zero otherwise, and

  o *lane ≥ 3* was equal to '1' if the pavement sample had three or more lanes and zero otherwise.

- NDOT classifies pavement samples under

  o The Interstate Route (IR),

  o The National Highway System (NHS), or

  o The Surface Transportation Program (STP).

  Two dummy variables, *nhs* and *stp*, were encoded as: *nhs* was equal to '1' if a segment belonged to the NHS; otherwise, *nhs* was equal to '0'. Similarly, *stp* was equal to '1' if a segment belonged to STP; otherwise, *stp* was equal to '0'.

- NDOT grouped its roadway network into five prioritization categories, 1- 5, using such factors as facility type and traffic volumes (NDOT, 2011). The type and

394            frequency of maintenance and rehabilitation (M&R) activities vary among these

395            prioritization categories. Four dummy variables – *category=2*, *category=3*,

396            *category=4*, and *category=5* – were encoded as:

397                   o   *category=2* was equal to '1' if the pavement sample is under Prioritization

398                      Category 2, '0' otherwise; and

399                   o   *category=3* was equal to '1' if the pavement sample was under

400                      Prioritization Category 3, '0' otherwise.

401            The same approach was used for the other three dummy variables.

402        •   Code of Federal Regulations (CFR) Title 23 part 470 mandates state agencies to

403            identify the functional class of roads and streets. NDOT divided its roadway network

404            into seven functional classes: (i) Interstate and Highway, (ii) Other Freeways and

405            Expressway, (iii) Principal Arterial-Other, (iv) Minor Arterial, (v) Major Collector,

406            (vi) Minor Collector, and (vii) Local. Six dummy variables – *f_class=2*, *f_class=3*,

407            *f_class=4, f_class=5, f_class=6* and *f_class=7* – were encoded as follows:

408                   o   *f_class=2* was equal to '1' if the pavement sample was an Interstate and

409                      Highway, '0' otherwise;

410                   o   *f_class=3* was equal to '1' if the pavement sample was classified as Other

411                      Freeways and Expressway, '0' otherwise.

412            The same approach was used for the other four dummy variables.

413       A total of 4,138 samples – including 14,638 observations from 2001 to 2010 and 3,005

414 observations from 2011 and 2012 – were available for model estimation and validation,

415 respectively. Table 1 illustrates a subset of data used in the experiments.

416    **Estimation parameters**

417    The existing literature does not provide hard-and-fast rules to define the limiting VIF beyond the

418    one that indicates a serious multicollinearity problem (Petraitis et al. 1996). Many studies (Myers

419    1990; Neter et al. 1996; Chatterjee and Hadi 2000) suggested that a multicollinearity problem

420    was serious if the VIF was greater than 10. In this study, all explanatory variables with VIF > 10

421    were excluded from the final models. Other estimation parameters that were required were set by

422    using previous experience (Paz et al. 2015a; Paz et al. 2015b) and sensitivity analyses. Table 2

423    provides the parameter values used in this study.

424    **Experiment results and discussion**

425    Function F in the constraint expressed by (13) was used to determine the maximum number of

426    feasible clusters for the dataset used in this study. The algorithm found that 16 was the maximum

427    number of feasible clusters that fulfilled the requirements imposed by the constraints for feasible

428    partitions.

429         The solution algorithm proposed in the section, Solution to the Mathematical Program,

430    sought for the optimum number of clusters by exploring each of all feasible clusters (i.e., $K = 2$

431    to 16). Thus, the algorithm determined that 6-cluster CLR models provided the optimum solution

432    with the lowest BIC. Figure 3a shows the BIC trend over the number of clusters that were

433    considered in this experiment. Figure 3b shows the convergence of the objective function, BIC,

434    over iterations when the six-cluster CLR models were used. After 983 iterations, the BIC

435    decreased from the initial value of 9,283 to the final value of 6,443, with an improvement of

436    31%.

437         Coefficients for the variables, *trucks* and *freeze_thaw*, were positive. This is counter-

438    intuitive because a pavement deteriorates faster when it is subjected to heavy trucks and frequent

439     freeze-thaw cycles. Hence, additional data analysis was performed to investigate the data quality.

440     The analysis showed average positive trends of PSI for these variables; this could be because

441     pavements having a larger number of trucks and freeze-thaw cycles often are designed to have

442     stronger pavement structures, and are continuously maintained. This research did not use any

443     explanatory variable that relates PSI to pavement structure. Hence, *trucks* and *freeze_thaw,*

444     which were likely to be positively correlated with missing information, such as pavement

445     structure, may have captured this hidden effect. There could be other reasons for the positive

446     coefficients for *trucks* and *freeze-thaw*; however, this investigation did not find enough evidence

447     to justify these positive trends. Hence, these two variables were excluded from the models, and

448     new model parameters were estimated. The effect of these two variables was discussed in

449     Khadka and Paz (2017a). Future research is recommended to investigate this issue. Table 3

450     provides the estimated parameters for 6-cluster models.

451         This study used a 5% significance level. Results showed that seven explanatory variables

452     – *elevation*, *precip*, *min_temp*, *max_temp*, *wet_days*, *nhs*, and *stp* – were not included in any of

453     the final estimated 6-cluster models. As the constraints for significant variables were imposed,

454     the algorithm excluded these seven variables because they were either associated with high VIF,

455     causing multicollinearity, or were statistically insignificant. Hence, the resultant models only had

456     statistically significant explanatory variables. Table 4 shows the binary matrix, *V*, associated with

457     the 6-cluster models estimated in this study. Each '1' indicates a 'significant' variable for a

458     particular cluster; '0' indicates otherwise.

459         Table 3 also includes the VIFs of the significant explanatory variables. All the VIF values

460     were less than five, which indicated that the associated explanatory variables in each model did

461    not have strong correlations among each other. Hence, the resultant models were free from

462    serious multicollinearity problems.

463    The six models included different significant explanatory variables. In addition, the

464    common variables had different estimated coefficients. These differences indicated that

465    pavement samples across the clusters were heterogeneous by the effect of explanatory variables,

466    and exhibited different types of performance behavior. For example, the samples exhibited

467    different deterioration rates as they got older. The estimated coefficients for *age* were -0.039 and

468    -0.022 for Clusters #1 and #2, respectively. However, pavement samples in Clusters #1 and #2

469    performed similarly with respect to traffic-loading conditions. That is, the estimated coefficients

470    for *adt* in Clusters #1 and #2 were -0.013 and -0.012, respectively.

471    Only four variables – *intercept*, *age*, *adt*, and *rut_depth* – were common for all six

472    models; and all of them had a negative sign, except for the intercept. All the estimated intercept

473    values were realistic. The PSI of a newly constructed pavement was about 4.5 (Christopher et al.

474    2006). However, the intercepts differed across the models. The negative signs of *age* and *adt*

475    indicated that the conditions deteriorated when a pavement became older and was subjected to

476    greater traffic loadings, respectively. Similarly, the PSI of a pavement sample decreased as

477    rutting along the pavement became deeper.

478    It was observed that Clusters #2 to #5, which had as significant variables *category=2*,

479    *category=3*, *category=4*, and *category=5*, also had variables *lane≤2* and *lane≥3* as significant. In

480    contrast, the variable *f_class* was not significant in these clusters. The estimated coefficients of

481    the variables *category=2*, *category=3*, *category=4*, and *category=5* were negative, and the

482    coefficient increased as the category level went up. This indicated that the average PSIs in these

483    four category levels (i.e., from 2 to 5) were smaller than for that of Category 1, and decreased as

484    the level went up. This was expected, because NDOT assigned the highest priority – in terms of

485    maintaining good conditions – to the roadway segments identified as Category 1 and the lowest

486    priority to the roadway segments identified as Category 5 (NDOT 2011). The variable *f_class*

487    was significant only in Clusters 1 and 6.  The coefficients for all six classes were negative,

488    except for the *f_class=2* in Cluster 6. A positive sign indicated that the pavement samples

489    classified as Class 2 had a higher average PSI than for the segments classified as Class 1. It also

490    was observed that for both clusters, the coefficient increased as the class number went up, except

491    for the *f_class=7*. A possible reason was that the estimation was based on only 44 observations

492    (Functional Class 7), which might not represent actual conditions.

493    **Model performance**

494    In CLR, minimizing overall SSE translates the maximization of variations in the dependent

495    variable explained by clustering process and regression models (Brusco et al. 2008). CLR does

496    not differentiate between the variations explained by the clustering process and variations

497    explained by regression models. Hence, in some cases, variation in the dependent variable could

498    be minimized by the clustering process even if variations explained by the regression models are

499    small. This creates a potential for overfitting the data in cases when regression relationships are

500    not strong.

501        Brusco et al. (2008) proposed a procedure to diagnose the presence of overfitting in the

502    resultant CLR models. Five different metrics were calculated for the optimum 6-cluster models,

503    and are included in Table 5. The results showed that the between-clusters sum of squares

504    (BCSS), which represent the variations explained by the clustering process, was equal to '4',

505    which is less than 1% of the total sum of squares (TSS). The sum of squares due to regression

506    (SSR) was equal to 1,130, which was 47% of the within-clusters sum of squares (WCSS). The

507   WCSS represents the sum of the variations across clusters, which is the sum of SSE and SSR.

508   This indicated that there was no overfitting, as most of the variations in PSI was explained by

509   within-cluster regressions. However, SSE accounted for 53% of the TSS, which indicated that

510   the resultant models had relatively high errors, possibly due to the nature of the data. In addition,

511   the estimated linear function might not have been the best to use in order to explain the pavement

512   deterioration.

513         The prediction accuracy of the models was evaluated by calculating the RMSE, the

514   normalized root-mean-square error (NRMSE), and the mean absolute error (MAE), using (14),

515   (15), and (16), respectively:

516   $$RMSE = \sqrt{\sum_i^{\eta}(y_{it} - \hat{y}_{it})^2 / \eta} \tag{14}$$

517   $$NRMSE = \frac{RMSE}{y_{max} - y_{min}} \tag{15}$$

518   $$MAE = \frac{1}{\eta}\sum_1^{\eta}|y_{it} - \hat{y}_{it}| \tag{16}$$

519   where $y_{it}$ = the observed PSI, $\hat{y}_{it}$ = the predicted PSI, $y_{max}$ = the maximum observed PSI, $y_{min}$ =

520   the minimum observed PSI, and $\eta$ = the number of predictions. The estimated model coefficients

521   were applied to the test dataset described in the Data section to estimate PSIs for 2011 and 2012.

522   The overall RMSE, NRMSE, and MAE values for all the models were 0.47, 0.17, and 0.36,

523   respectively. This indicated that the resultant models were robust.

524         In addition, to diagnose the variation in the prediction errors, the RMSE, NRMSE, and

525   MAE were calculated separately for all six models. Table 6 provides the RMSE, NRMSE, and

526   MAE values for all the models as well as the individual models. It was observed that the

527   differences between RMSE and MAE values were approximately equal for all the models, which

528   indicated that the prediction errors were well distributed among the clusters.

529    Figure 4a shows a scattered plot of predicted versus observed PSIs for 2011 and 2012.

530    The degree of prediction error of the models was reflected by the relative positions of the data

531    points from the $45^0$ line. Data points above the $45^0$ line were over-predicted, while those under

532    the $45^0$ line were under-predicted. Results indicated that the predicted PSIs ranged from 2.70 to

533    4.42, whereas the observed PSIs ranged from 1.64 to about 4.44. In particular, the CLR models

534    overestimated PSIs that were at the lower end of the data. Possible reasons for overestimation

535    could be that this study did not include any explanatory variables that captured the pavement

536    structure. In addition, improvements by routine maintenance activities were ignored.

537    Figure 4b provides the percentages of observations that were within different ranges of

538    error. For example, about 74% of the total number of predictions were contained within a ±15%

539    range of error. Figure 5 shows individual scattered plots of predicted versus observed PSIs for all

540    six models.


541    **CONCLUSIONS AND RECOMMENDATIONS**

542    In this paper, a comprehensive mathematical program is proposed to estimate PPMs that

543    minimize the estimation error by simultaneously finding 1) the optimum number of pavement

544    clusters, 2) cluster memberships of the samples, 3) cluster-specific significant explanatory

545    variables, and 4) regression coefficients. To solve the mathematical program, Simulated

546    Annealing integrated with All-subset regression was implemented. The algorithm has the

547    capability to identify potential explanatory variables that cause serious multicollinearity in a

548    model.

549    VIF was used to measure the effect of multicollinearity in a model. In this study,

550    multicollinearity was addressed using a traditional approach where correlated variables were

551    removed one at a time until the effect of multicollinearity became minimal. However, a better

552     way to address multicollinearity is to consider the trade-off between removing and keeping

553     potential explanatory variables that are expected to cause multicollinearity. Future research is

554     recommended to integrate such an experiment in the CLR framework.

555     After addressing the multicollinearity issue, the proposed algorithm identified the

556     relevant explanatory variables to be included in the models. All possible combinations of the

557     explanatory variables were evaluated to select the best model for each cluster. Hence, the

558     estimated CLR models included cluster-specific significant explanatory variables that were free

559     from multicollinearity.

560     The algorithm explored all the feasible clusters that could be formed for the data used in

561     the experiments, and found that 6-cluster models were the optimum solution. The algorithm

562     determined the significant explanatory variables to be traffic-loading conditions of both ADT and

563     the number of trucks, age, rut-depth, function class, prioritization category, freeze-and-thaw

564     cycles, and the number of lanes. In the literature, all these variables were considered to be the

565     most critical factors for pavement deterioration (Saraf and Majidzadeh 1992; Prozzi and Madanat

566     2004; Kim and Kim 2006; Salama et al. 2006). Both the magnitude and sign of the estimated

567     regression coefficients were as expected, and were realistic. This indicates that the proposed

568     algorithm was very effective when selecting the explanatory variables that were relevant.

569     The estimated CLR models first were analyzed to investigate the presence of overfitting,

570     and the results showed that the models did not possess any overfitting issues. To investigate the

571     predictive capability of the models, RMSE, NRMSE, and MAE were calculated for all the

572     models as well as for individual models. The overall RMSE, NRMSE, and MAE values of 0.47,

573     0.17, and 0.36, respectively, indicated that the estimated models had small estimation errors. In

574     addition, the results showed that both the differences between the RMSE and MAE values for all

575     six models were approximately equal, which indicated that the prediction error was well

576     distributed among the models. Even so, the models still were associated with prediction errors.

577     The linear functional form used in this study did not exactly fit the data used in the

578     experiments. Hence, it would be worth investigating the proposed methodology by using

579     nonlinear relationships between pavement performance measures and multiple explanatory

580     variables. Various forms of power and sigmoidal models (Sadek et al. 1996; Luo and Chou 2006;

581     Zhang and Durango-Cohen 2014; and Chen and Mastin 2015) could be investigated.

582     Finally, the results indicated that each cluster had almost an equal number of members

583     (i.e., pavement samples). However, it is unlikely that the underlying clusters had equally

584     distributed pavement samples. An interesting aspect worthy of investigation would be to explore

585     the likelihood of distribution of the pavement samples and the associated physical characteristics.

586     Further investigation would be required to identify any dominant explanatory variables of the

587     pavement samples that define a cluster.

588     This study also proposed a heuristic to assign cluster membership to a pavement sample.

589     It was assumed that a new sample had observations for all explanatory variables included in the

590     estimated CLR models as well as the dependent variable (i.e. PSI). Future research is

591     recommended to develop a procedure to assign cluster membership to pavement samples that

592     were not included during model estimation and lack data about the dependent variable.

593     The proposed algorithm was designed to search for a global minimum; however, a large

594     amount of computational time is required. Another avenue for future research would be to

595     develop faster and more efficient combinatorial algorithms that would guarantee global

596     optimality.

604 **NOTATION**

605 *The following symbols are used in this paper:*

606 $I$ = Number of pavement samples in the network;

607 $i$ = Subscript for a pavement sample in the network, $i \in I$;

608 $T_i$ = Number of observation periods for a pavement sample $i$;

609 $t$ = Subscript for an observation period for a pavement sample $i$, $t \in T_i$;

610 $O$ = Total number of observations = $\sum_i^I T_i \ \forall \ i \in I$;

611 $J$ = Number of explanatory variables;

612 $j$ = Subscript for an explanatory variable including an intercept, $j = 0, \dots, J$

613 $x_{ijt}$ = Measurement of an explanatory variable $j$ for a sample $i$ at observation period $t$ that is

614       assigned to a cluster $k \ \forall \ i \in I, j \in J, t \in T_i$;

615 $y_{it}$ = Measurement of dependent variable for a sample $i$ at observation period $t$ that is assigned to

616       a cluster $k \ \forall \ i \in I, t \in T_i$;

617 $K$ = Optimum number of clusters ($1 \leq k \leq K_{max}$);

618 $k$ = Subscript for a cluster, $k \in K$;

619 $K_{max}$ = Maximum number of potential clusters that could be formed using the given data;

620 $n$ = Minimum number of observations required in a cluster;

621 $C_k$ = Set of pavement samples that are assigned to cluster $k \ \forall \ k \in K$;

622 $\delta$ = Total number of significant explanatory variables including intercepts in all clusters;

623 $v_{jk}$ = Binary indicator that represents significance of an explanatory variable including an

624 intercept in a cluster $k \ \forall \ j = 0, \dots, J, k \in K$;

625 $p_{ik}$ = Cluster membership of a pavement sample $i$ to a cluster $k$, $\forall \ i \in I, k \in K$;

626 $\beta_{jk}$ = Estimated regression coefficient for an explanatory variable $j$ including an intercept in

627 cluster $k \ \forall \ j = 0, \dots, J, k \in K$;

628 **REFERENCES**

629 Aarts, E., Korst, J. and Michiels, W. (2005). "Simulated Annealing." *Search Methodologies:*

630 *Introductory Tutorials in Optimization and Decision Support Techniques*, Springer US,

631 187–210, DOI: 10.1007/0-387-28356-0_7.

632 Abdul-Wahaba, S. A., Bakheitb, C. S. and Al-Alawia, S. M. (2005). "Principal component and

633 multiple regression analysis in modelling of ground-level ozone and factors affecting its

634 concentrations." *Environmental Modelling & Software*, 20(10), 1263–1271.

635 Attoh-Okine, N. and Adarkwa, O. (2013). "Pavement Condition Surveys–Overview of Current

636 Practices." 〈www.sites.udel.edu/dct/files/2013/10/Rpt-245-Pavement-Condition-Okine-

637 DCTR422232-1pzk0uz.pdf〉 (January 2016).

638 Baumann, K. (2003). "Cross-validation as the objective function for variable-selection

639 techniques." *TrAC Trends in Analytical Chemistry*, 22(6), 395–406, DOI:

640 10.1016/S0165-9936(03)00607-1.

641 Berk, K. N. (1978). "Comparing Subset Regression Procedures." *Technometrics*, 20(1), 1–6.

642 Brusco, M. J. (2014). "A comparison of simulated annealing algorithms for variable selection in

643 principal component analysis and discriminant analysis." *Computational Statistics &*

644 *Data Analysis*, 77, 38–53.

645 Brusco, M. J., Cradit, J. D., Steinley, D. and Fox, G. L. (2008). "Cautionary Remarks on the Use

646 of Clusterwise Regression." *Multivariate Behavioral Research*, 43(1), 29–49.

647 Carbonnea, R. A., Caporossi, G. and Hansen, P. (2011). "Globally optimal clusterwise regression

648 by mixed logical-quadratic programming." *European Journal of Operational Research*,

649 212(1), 213–222.

650 Cěrny, V. (1985). "A Thermodynamical Approach to the Travelling Salesman Problem; An

651 Efficient Simulation Algorithm." *J. of Optimization Theory and Applic.* 45, 41-55.

652 Chatterjee, S. and Hadi, A. S. (2000). *Regression analysis by example*, John Wiley and Sons,

653 New York, USA.

654 Chen, D. and Mastin, N. (2015). "Sigmoidal models for predicting pavement performance

655 conditions." *J. Perform. Constr. Facil.*, 30(4). DOI:10.1061/(ASCE)CF.1943-

656 5509.0000833.

657 Christopher, B. R., Schwartz, C. and Boudreau, R. (2006). "Geotechnical Aspects of

658 Pavements." *Report No. FHWA NHI-05-037,* U.S. Department of Transportation,

659 Washington, D.C.

660 Davies, R. M. and Sorenson, J. (2000). "Pavement preservation: Preserving our investment in

661 highways." *Public Roads*, 63(2), 63–69.

662 DeSarbo, W. S. and Corn, W. L. (1988). "A Maximum Likelihood Methodology for Clusterwise

663 Linear Regression." *Journal of Classification*, 5(2), 249–282.

664 DeSarbo, W. S., Oliver, R. L. and Rangaswamy, A. (1989). "A Simulated Annealing
665     Methodology for Clusterwise Linear Regression." *Psychometrika*, 54(4), 707–736.

666 Dolan, W. B., Cummings, P. T. and LeVan, M. D. (1989). "Process optimization via simulated
667     annealing: Application to network design." *AIChE Journal*, 35(5), 725–736, DOI:
668     10.1002/aic.690350504.

669 Dormann, C. F., Elith, J., Bacher, S., Buchmann, C., Carl, G., Carré, G., Marquéz, J. R. G.,
670     Gruber, B., Lafourcade, B., Leitão, P. J. and Münkemüller, T. (2013). "Collinearity: a
671     review of methods to deal with it and a simulation study evaluating their performance."
672     *Ecography*, 36(1), 27-46.

673 Efron, B., Hastie, T., Johnstone, I. and Tibshirani, R. (2004). "Least Angle Regression." *The
674     Annals of Statistics*, 32(2), 407–499.

675 Fowlkes, E. B., Gnanadesikan, R. and Kettenring, J. R. (1988). "Variable Selection in
676     Clustering." *Journal of Classification*, 5(2), 205–228.

677 Garside, M. J. (1965). "The Best Subset in Multiple Regression Analysis." *Applied Stat. Journ.
678     of the Royal Statist. Society, Series C.*, 14(2), 196–200.

679 Geweke, J. F., and Meese, R. A. (1981). "Estimating Regression Models of Finite but Unknown
680     Order." *International Economics Review*, 22, 55–70.

681 Gorman, J. W. and Toman, R. J. (1966). "Selection of Variables for Fitting Equations to Data."
682     *Technometrics*, 8(1), 27–51, DOI: 0.1080/00401706.1966.10490322.

683 Gunst, R. F. and Webster, J. T. (1975). "Regression analysis and problems of multicollinearity."
684     *Communications in Statistics*, 4(3), 277–292, DOI: 10.1080/03610927308827246.

685     Gupta, M. and Ibrahim, J. G. (2007). "Variable selection in regression mixture modeling for the

686          discovery of gene regulatory networks." *Journal of the American Statistical Association*,

687          102(479), 867–880.

688     Hajj, E. Y., Loria, L. and Sebaaly, P. E. (2010). "Performance Evaluation of Asphalt Pavement

689          Preservation Activities." *Transportation Research Record: Journal of the Transportation*

690          *Research Board*, 2150, 36–46.

691     Harrell, F. E. (2001). *Regression Modeling Strategies: with Applications to Linear Models,*

692          *Logistic Regression, and Survival Analysis (Second Ed.)*, Springer, Cham, Heidelberg,

693          New York Dordrecht, London.

694     Highway Research Board. (1962). "The AASHO road test." *Special Rep. No. 61A-E*, National

695          Academy of Science, National Research Council, Washington, DC.

696     Hocking, R. R. and Leslie, R. N. (1967). "Selection of the Best Subset in Regression Analysis."

697          *Technometrics*, 9(4), 531–540.

698     Hudson, W. R., Haas, R. and Perrone, E. (2015). "Measures of Pavement Performance must

699          consider the Road User." *Proc.,9th International Conference on Managing Pavement*

700          *Assets*, Alexandria, VA.

701     Johnson, D. S., Aragon, C. R., McGeoch, L. A. and Schevon, C. (1989). "Optimization by

702          simulated annealing: an experimental evaluation; part I, graph partitioning." *Operations*

703          *research*, 37(6), 865–892.

704     Joseph B. Kadane, J. B. and Lazar, N. A. (2004). "Methods and Criteria for Model Selection."

705          *Journal of the American Statistical Association*, 99(465) 279-290.

706     Katz, R. W. (1981). "On Some Criteria for Estimating the Order of a Markov Chain."

707          *Technometrics*, 23, 243–249.

708  Ketchen, D. J. and Shook, C. L. (1996). "The Application of Cluster Analysis in Strategic
709       Management Research: An Analysis and Critique." *Strategic Management Journal*,
710       17(6), 441–458.

711  Khadka, M. and Paz, A. (2017a). "Limitations of Existing Pavement Performance Models and a
712       Potential Solution." *World Conference on Pavement and Asset Management*, Baveno,
713       Italy.

714  Khadka, M., & Paz, A. (2017b). Estimation of optimal pavement performance models for
715       highways. Tenth International Conference on the Bearing Capacity of Roads, Railways
716       and Airfields. Athens, Greece.

717  Khadka, M., Paz, A., Cristian, A., and Hale, D. K. (2017). "Simultaneous Generation of
718       Optimum Pavement Clusters and Associated Performance Models." Manuscript
719       submitted for publication in *Transportmetrica A: Transport Science*. Under review.

720  Kim, S. and Kim N. (2006). "Development of performance prediction models in flexible
721       pavement using regression analysis method." *KSCE J. of Civil Engineering*, 10(2), 91–
722       96.

723  Kim, Y., Kwon, S., and Choi, H. (2012). "Consistent Model Selection Criteria on High
724       Dimensions." *Journal of Machine Learning Research*, 13(1), 1037-1057.

725  Kirkpatrick, S., Gelatt, C. D. and Vecchi, M. P. (1983). "Optimization by Simulated Annealing."
726       *Science, New Series,* 220(4598), 671-680.

727  Koehler, A. B. and Murphree, E. S. (1988). "A Comparison of the Akaike and Schwarz Criteria
728       for Selecting Model Order." *Applied Statistics*, 37, 187-19.

729  Labi, S. and Sinha, K. C. (2003). "Life-cycle evaluation of highway pavement preventive
730       maintenance." [CD-ROM], *82$^{nd}$ annual meeting of the Transportation Research Board*,

National Research Council, Washington DC.

Lau, K., Leung, P. and Tse, K. (1999). "A mathematical programming approach to clusterwise regression model and its extensions." *European Journal of Operational Research*, 116(3), 640–652.

Liu, H. H. and Ong, C. S. (2008). "Variable selection in clustering for marketing segmentation using genetic algorithms." *Expert Systems with Applications*, 34(1), 502–510.

Lu, H., Huang, S., Li, Y. and Yang, Y. (2014). "Panel Data Analysis Via Variable Selection and Subject Clustering." *Data Mining for Services*, 3, Springer-Verlag, 31–76.

Luo, Z. and Chou, E. Y. J. (2006). "Pavement Condition Prediction Using Clusterwise Regression." *Transportation Research Record: Journal of the Transportation Research Board*, No. 1974, 70–77.

Luo, Z. and Yin, H. (2008). Probabilistic Analysis of Pavement Distress Ratings with the Clusterwise Regression Method, *Transportation Research Record: Journal of the Transportation Research Board*, No. 2084, 38–46.

Mallows, C. L. (1973). "Some Comments on Cp." *Technometrics*, 15(4), 661–675.

Mansfield, E. R. and Helms, B. P. (1982). "Detecting Multicollinearity." *The American Statistician*, 36(3a), 158-160. DOI: 10.1080/00031305.1982.10482818.

Marquardt, D. W. (1970). "Generalized inverses, ridge regression, biased linear estimation, and nonlinear estimation." *Technometrics*, 12(3), 591–612.

Maugis, M., Celeux, G. and Martin-Magniette, M. L. (2009). "Variable Selection for Clustering with Gaussian Mixture Models." *Journal of the International biometric society*, 65(2), 701-709.

753 Maydeu-Olivares, A. and García-Forero, C. (2010). "Goodness-of-fit testing." *International*
754 *Encyclopedia of Education*, 7, 190–196.

755 Mehmood, T., Liland, K. H., Snipen, L. and Sæbø, S. (2012). "A review of variable selection
756 methods in Partial Least Squares Regression." *Chemometrics and Intelligent Laboratory*
757 *Systems*, 118, 62–69.

758 Midi, H., Sarkar, S. K. and Rana, S. (2010). "Collinearity diagnostics of binary logistic
759 regression model." *Journal of Interdisciplinary Mathematics*, 13(3), 253–267.

760 Myers, R.H. (1990). *Classical and Modern Regression with Applications, 2nd edition*, PWS
761 Kent, Boston, MA.

762 Neter, J., Kutner, M. H., Nachtsheim, C. J. and Wasserman, W. (1996). *Applied linear statistical*
763 *models*, Irwin, Chicago, Illinois, USA.

764 Nevada Department of Transportation (2011). *Pavement Management System Overview*,
765 Material Division, Nevada Department of Transportation, USA.

766 Nikolaev A. and Jacobson S. (2010). "Simulated annealing." *Handbook of Metaheuristics,*
767 *International Series in Operations Research and Management Science*, 146, Springer,
768 Berlin, 1–39, DOI: 10.1007/978-1-4419-1665-5_1.

769 Ohlemuller, R., Anderson, B. J., Araujo, M. B., Butchart, S. H., Kudrna, O., Ridgely, R. S., and
770 Thomas. C. D. (2008). "The coincidence of climactic and species rarity: high risk to
771 small range species from climate change." *Biology Letters*, 4, 568–572.

772 Park, Y. W., Jiang, Y., Klabjan, D. and Williams, L. (2015). "Algorithms for Generalized
773 Cluster-wise Linear Regression." 〈http://www.dynresmanagement.com/publications〉
774 (January 2016).

775    Paz, A., Molano, V. and Sanchez, M. (2015a). "Holistic Calibration of Microscopic Traffic

776        Flow Models: Methodology and Real World Application Studies." *Engineering and*

777        *Applied Sciences Optimization: Dedicated to the memory of Professor M.G. Karlaftis*,

778        38, Ed. 1. Springer International Publishing.

779    Paz, A., Molano, V., Martinez, E., Gaviria, C. and Arteaga, C. (2015b). "Calibration of Traffic

780        Flow Models Using a Memetic Algorithm." *Transportation Research Part-C: Emerging*

781        *Technologies*, 55, 432–443.

782    Petraitis, P. S., Dunham, A. E. and Niewiarowski, P. H. (1996). "Inferring multiple causality: the

783        limitations of path analysis." *Functional Ecology*, 10(4), 421–431.

784    Prozzi, J. A. and Madanat, S. M. (2004). "Development of pavement performance models by

785        combining experimental and field data." *J. of Infrastructure Systems*, 10(1), 9–22.

786    Rao, C. R. and Wu, Y. (1989). "A strongly consistent procedure for model selection in a

787        regression problem." *Biometrika*, 76(2), 369–374.

788    Román-Román, P., Romero, D., Rubio, M. A. and Torres-Ruiz, F. (2012). "Estimating the

789        parameters of a Gompertz-type diffusion process by means of Simulated Annealing."

790        *Applied Mathematics and Computation*, 218(9), 5121–5131.

791    Rutenbar, R.A. (1989). "Simulated Annealing Algorithms: An Overview." *IEEE Circuits and*

792        *Devices Magazine*, 5(1), 19–26.

793    Sadek, A. W., Freeman, T. E. and Demetsky, M. J. (1996). "Deterioration prediction modeling of

794        Virginia's interstate highway system." *Transportation Research Record: Journal of the*

795        *Transportation Research Board*, No. 1524, 118–129.

796    Salama, H., Chatti, K. and Lyles, R. (2006). "Effect of Heavy Multiple Axle Trucks on Flexible

797         Pavement Damage Using In-Service Pavement Performance Data." *J. Transp. Eng.*,

798         132(10), 763-770.

799    Saraf C. L. and Majizzadeh, K. (1992). "Distress prediction models for a network level PMS."

800         *Transportation Research Record: Journal of the Transportation Research Board*, No.

801         1344, 38–48.

802    Schlittgen, R. (2011). "A weighted least-squares approach to clusterwise regression." *Advances*

803         *Statistical Analysis*, 95(2), 205–217.

804    Schwarz, G. (1978). "Estimating the dimension of a model." *Annals of Statistics*, 6(2), 461–464.

805    Shoukry, S. N., Martinelli, D. R. and Reigle, J. A. (1997). "Universal Pavement Distress

806         Evaluator Based on Fuzzy Sets." *Transportation Research Record: Journal of the*

807         *Transportation Research Board*, No. 1592, 180–186.

808    Spath, H. (1979). "Algorithm 39: Clusterwise linear regression." *Computing*, 22(4), 367–373.

809    Sridhar, J. and Rajendran, C. (1993). "Scheduling in a cellular manufacturing system: a

810         simulated annealing approach." *International Journal of Production Research*, 31(12),

811         2927-2945, DOI: 10.1080/00207549308956908

812    Tacq. J. (1997). *Multivariate analysis techniques in social science research: From problem to*

813         *analysis*, Sage Publications, London.

814    Tan, T., Suk, H. W., Hwang, H. and Lim, J. (2013). "Functional fuzzy clusterwise regression

815         analysis." *Adv. Data Anal. Classif.*, 7(1), 57–82.

816    Terzi, S. (2006). "Modeling the Pavement Present Serviceability Index of Flexible Highway

817         Pavements Using Data Mining." *Journal of Applied Sciences*, 6(1), 193–197.

818    Thompson, M. L. (1978). "Selection of Variables in Multiple Regression: Part I. A Review and

819        Evaluation." *International Statistical Review*, 46(1), 1–19, DOI: 10.2307/1402505.

820    Tibshirani, R. (1996). "Regression shrinkage and selection via the lasso." *J. R. Statist. Soc. B*, 58,

821        267–288.

822    Uzoma, U. and Jeremiah, O. O. (2016). "An Alternative Approach to AIC and Mallow's Cp

823        Statistic-Based Relative Influence Measures (RIMS) in Regression Variable Selection."

824        *Open Journal of Statistics*, 6(1), 70-75, DOI: 10.4236/ojs.2016.61009

825    Vrieze, S. I. (2012). "Model selection and psychological theory: A discussion of the differences

826        between the Akaike Information Criterion (AIC) and the Bayesian Information Criterion

827        (BIC)." *Psychol Methods*, 17(2), 228–243.

828    Washington, S. P., Karlaftis, M. G., and Mannering, F. L. (2011). *Statistical and Econometric*

829        *Methods for Transportation Data Analysis (Second Ed.)*, Chapman Hall/CRC, Boca

830        Raton, FL.

831    Wedel, M. and SteenKamp, J. (1989). "Fuzzy clusterwise regression approach to benefit

832        segmentation." *Int. J. Res. Mark.*, 6(4), 241–258.

833    Yang, Y. (2005). "Can the strengths of AIC and BIC be shared? A conflict between model

834        identification and regression estimation." *Biometrika*, 92(4), 937–950.

835    Yoo, W., Mayberry, R., Bae, S., Singh, K., (Peter) He, Q. and Lillard, J. W. (2014). "A Study of

836        Effects of MultiCollinearity in the Multivariable Analysis." *Int J Appl Sci Technol.*, 4(5),

837        9–19.

838    Zhang, W. and Durango-Cohen, P. (2014). "Explaining Heterogeneity in Pavement

839        Deterioration: Clusterwise Linear Regression Model." *J. Infrastruct. Syst.*, 20(2). DOI:

840        10.1061/(ASCE)IS.1943-555X.0000182.

841 Zhen, Z., Yan, L. and Nan, K. (2012). "Clusterwise linear regression with the least sum of

842   absolute deviations - An MIP approach." *Int. J. Oper. Res.*, 9(3), 62172.

843

844 **Table 1.** A Subset of Data used in the Experiments

| Sample ID | 1 | | | | 2 | | | |
|---|---|---|---|---|---|---|---|---|
| Year | 2001 | 2002 | 2003 | 2004 | 2001 | 2002 | 2003 | 2004 |
| psi | 3.82 | 3.73 | 3.71 | 3.62 | 3.96 | 3.88 | 3.86 | 3.53 |
| age | 0 | 1 | 2 | 3 | 0 | 1 | 2 | 3 |
| adt | 725 | 825 | 950 | 950 | 725 | 825 | 950 | 950 |
| trucks | 20 | 20 | 19 | 20 | 20 | 20 | 19 | 20 |
| elevation | 4750 | 4750 | 4750 | 4750 | 4750 | 4750 | 4750 | 4750 |
| precip | 8.25 | 8.25 | 6.65 | 6.65 | 6.65 | 6.65 | 6.65 | 6.65 |
| min_temp | 33 | 33 | 36 | 36 | 36 | 36 | 36 | 36 |
| max_temp | 65 | 65 | 67 | 67 | 67 | 67 | 67 | 67 |
| wet_days | 45 | 45 | 41 | 41 | 41 | 41 | 41 | 41 |
| freeze_thaw | 176 | 176 | 154 | 154 | 154 | 154 | 154 | 154 |
| rut_depth | 0.09 | 0.08 | 0.08 | 0.09 | 0.01 | 0.01 | 0.01 | 0.01 |
| lane≤2 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| lane≥3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| nhs | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| stp | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| f_class=2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| f_class=3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| f_class=4 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| f_class=5 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| f_class=6 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| f_class=7 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Category=2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Category=3 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Category=4 | 0 | 1 | 1 | 1 | 0 | 1 | 1 | 1 |
| Category=5 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

845

846     **Table 2.** Estimation Parameters Used in the Experiments

| Parameter | Value | Remarks |
|---|---|---|
| $\theta_0$ | 10 | Initial temperature |
| $\theta_{min}$ | 10e-17 | Final minimum temperature |
| $B$ | 30 | Boltzmann constant |
| $\lambda$ | 0.97 | Cooling rate |
| $N_{max}$ | 5 | Number of neighborhood solutions generated at each temperature level |
| $n$ | 800 | Minimum number of observations required in a cluster |
| $N_{ps}$ | 100 | Number of pavement samples, which memberships were changed to generate a neighborhood cluster |
| $VIF_{max}$ | 10 | Limiting VIF |
| $\alpha$ | 5% | Level of Significance |

847

848    **Table 3.** Estimated Model Parameters using the Proposed CLR Approach

| Parameters | Cluster #1 | | | Cluster #2 | | | Cluster #3 | | |
|---|---|---|---|---|---|---|---|---|---|
| | $\beta_{j1}$ | VIF | p-value | $\beta_{j2}$ | VIF | p-value | $\beta_{j3}$ | VIF | p-value |
| *intercept* | 4.392 | - | < 0.0001 | 4.552 | - | < 0.0001 | 4.674 | - | < 0.0001 |
| *age* | -0.039 | 1.0 | < 0.0001 | -0.022 | 1.0 | < 0.0001 | -0.028 | 1.0 | < 0.0001 |
| *adt†* | -0.013 | 1.2 | < 0.0001 | -0.012 | 1.8 | < 0.0001 | -0.008 | 2.2 | < 0.0001 |
| *rut_depth* | -1.293 | 1.1 | < 0.0001 | -2.814 | 1.1 | < 0.0001 | -3.338 | 1.1 | < 0.0001 |
| *lane≤2* | - | - | - | -0.191 | 4.4 | < 0.0001 | -0.358 | 4.4 | < 0.0001 |
| *lane≥3* | - | - | - | -0.202 | 1.8 | < 0.0001 | -0.289 | 2.5 | < 0.0001 |
| *f_class=2* | -0.185 | 1.0 | 0.002 | - | - | - | - | - | - |
| *f_class=3* | -0.110 | 1.6 | < 0.0001 | - | - | - | - | - | - |
| *f_class=4* | -0.259 | 1.5 | < 0.0001 | - | - | - | - | - | - |
| *f_class=5* | -1.052 | 1.4 | < 0.0001 | - | - | - | - | - | - |
| *f_class=6* | -1.181 | 1.1 | < 0.0001 | - | - | - | - | - | - |
| *f_class=7* | -0.284 | 1.0 | 0.006 | - | - | - | - | - | - |
| *category=2* | - | - | - | -0.202 | 2.6 | < 0.0001 | -0.325 | 2.8 | < 0.0001 |
| *category=3* | - | - | - | -0.323 | 4.2 | < 0.0001 | -0.465 | 4.4 | < 0.0001 |
| *category=4* | - | - | - | -0.664 | 2.6 | < 0.0001 | -0.684 | 2.9 | < 0.0001 |
| *category=5* | - | - | - | -1.149 | 2.8 | < 0.0001 | -0.808 | 2.8 | < 0.0001 |
| *No. of Obs.* | 2,376 | | | 2,483 | | | 2,442 | | |
| *BIC* | 658 | | | 1,069 | | | 1,470 | | |

| Parameters | Cluster #4 | | | Cluster #5 | | | Cluster #6 | | |
|---|---|---|---|---|---|---|---|---|---|
| | $\beta_{j4}$ | VIF | p-value | $\beta_{j5}$ | VIF | p-value | $\beta_{j6}$ | VIF | p-value |
| *intercept* | 4.605 | - | < 0.0001 | 4.557 | - | < 0.0001 | 4.401 | - | < 0.0001 |
| *age* | -0.033 | 1.0 | < 0.0001 | -0.028 | 1.0 | < 0.0001 | -0.037 | 1.0 | < 0.0001 |
| *adt†* | -0.006 | 1.8 | < 0.0001 | -0.005 | 2.2 | < 0.0001 | -0.013 | 1.4 | < 0.0001 |
| *rut_depth* | -3.706 | 1.1 | < 0.0001 | -3.289 | 1.1 | < 0.0001 | -2.291 | 1.0 | < 0.0001 |
| *lane≤2* | -0.213 | 4.8 | < 0.0001 | -0.260 | 4.9 | < 0.0001 | - | - | - |
| *lane≥3* | -0.405 | 1.9 | < 0.0001 | -0.294 | 2.4 | < 0.0001 | - | - | - |
| *f_class=2* | - | - | - | - | - | - | 0.468 | 1.2 | < 0.0001 |
| *f_class=3* | - | - | - | - | - | - | -0.086 | 1.5 | < 0.0001 |
| *f_class=4* | - | - | - | - | - | - | -0.258 | 1.4 | < 0.0001 |
| *f_class=5* | - | - | - | - | - | - | -0.864 | 1.3 | < 0.0001 |
| *f_class=6* | - | - | - | - | - | - | -1.288 | 1.1 | < 0.0001 |
| *f_class=7* | - | - | - | - | - | - | -0.634 | 1.0 | < 0.0001 |
| *category=2* | -0.263 | 3.0 | < 0.0001 | -0.194 | 2.7 | < 0.0001 | - | - | - |
| *category=3* | -0.325 | 4.0 | < 0.0001 | -0.287 | 4.2 | < 0.0001 | - | - | - |
| *category=4* | -0.650 | 3.2 | < 0.0001 | -0.639 | 3.1 | < 0.0001 | - | - | - |
| *category=5* | -0.808 | 3.4 | < 0.0001 | -1.130 | 2.9 | < 0.0001 | - | - | - |
| *No. of Obs.* | 2,414 | | | 2,340 | | | 2,583 | | |
| *BIC* | 1,009 | | | 1,273 | | | 870 | | |

849    Note: *†* = variable value in thousands, and - = Not applicable

850

851   **Table 4.** Binary Matrix Showing Significance of the Explanatory Variables in the Estimated 6-

852   Cluster Models

| | Cluster | | | | | |
|---|---|---|---|---|---|---|
| Explanatory Variables | $k = 1$ | $k = 2$ | $k = 3$ | $k = 4$ | $k = 5$ | $k = 6$ |
| intercept | 1 | 1 | 1 | 1 | 1 | 1 |
| age | 1 | 1 | 1 | 1 | 1 | 1 |
| adt | 1 | 1 | 1 | 1 | 1 | 1 |
| elevation | 0 | 0 | 0 | 0 | 0 | 0 |
| precip | 0 | 0 | 0 | 0 | 0 | 0 |
| min_temp | 0 | 0 | 0 | 0 | 0 | 0 |
| max_temp | 0 | 0 | 0 | 0 | 0 | 0 |
| wet_days | 0 | 0 | 0 | 0 | 0 | 0 |
| rut_depth | 1 | 1 | 1 | 1 | 1 | 1 |
| lane≤2 | 0 | 1 | 1 | 1 | 1 | 0 |
| lane≥3 | 0 | 1 | 1 | 1 | 1 | 0 |
| nhs | 0 | 0 | 0 | 0 | 0 | 0 |
| stp | 0 | 0 | 0 | 0 | 0 | 0 |
| f_class=2 | 1 | 0 | 0 | 0 | 0 | 1 |
| f_class=3 | 1 | 0 | 0 | 0 | 0 | 1 |
| f_class=4 | 1 | 0 | 0 | 0 | 0 | 1 |
| f_class=5 | 1 | 0 | 0 | 0 | 0 | 1 |
| f_class=6 | 1 | 0 | 0 | 0 | 0 | 1 |
| f_class=7 | 1 | 0 | 0 | 0 | 0 | 1 |
| Category=2 | 0 | 1 | 1 | 1 | 1 | 0 |
| Category=3 | 0 | 1 | 1 | 1 | 1 | 0 |
| Category=4 | 0 | 1 | 1 | 1 | 1 | 0 |
| Category=5 | 0 | 1 | 1 | 1 | 1 | 0 |

853

854     **Table 5.** Metrics Calculated to Investigate the Presence of Overfitting in the Models

| Metric | Value | Remarks |
|--------|-------|---------|
| TSS | 2,419 | - |
| BCSS | 4 | 0.17% of TSS |
| WCSS | 2,415 | - |
| SSR | 1,130 | 47% of WCSS |
| SSE | 1,284 | 53% of TSS |

855

856    **Table 6.** RMSE, NRMSE, and MAE for Each Cluster

| Metric | Cluster | | | | | | Overall |
|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 | 6 | |
| RMSE | 0.47 | 0.46 | 0.49 | 0.47 | 0.48 | 0.49 | 0.47 |
| NRMSE | 0.18 | 0.18 | 0.18 | 0.17 | 0.18 | 0.19 | 0.17 |
| MAE | 0.37 | 0.37 | 0.37 | 0.35 | 0.36 | 0.38 | 0.36 |

857