

Copyright
by
Moo Yeon Kim
2016

**The Report Committee for Moo Yeon Kim
Certifies that this is the approved version of the following report:**

Segmentation of Highway Networks for Maintenance Operations

**APPROVED BY
SUPERVISING COMMITTEE:**

Supervisor:

Sinead Williamson

Jorge A. Prozzi

Segmentation of Highway Networks for Maintenance Operations

by

Moo Yeon Kim, B.E.; M.S.E.

Report

Presented to the Faculty of the Graduate School of

The University of Texas at Austin

in Partial Fulfillment

of the Requirements

for the Degree of

Master of Science in Statistics

The University of Texas at Austin

May 2016

Dedication

To my lovely wife Inae and my son Ian.

Acknowledgements

I would like to acknowledge my supervisor, Dr. Sinead Williamson, for her guidance, patient, and time. Without her encouragement I'd never have thought of pursuing the second master's degree. I would also like to acknowledge my reader, Dr. Jorge A. Prozzi, for his unfailing support as an academic supervisor.

Abstract

Segmentation of Highway Networks for Maintenance Operations

Moo Yeon Kim, M.S.Stat.

The University of Texas at Austin, 2016

Supervisor: Sinead Williamson

Pavement maintenance and rehabilitation (M&R) is important for transportation agencies to have a sustainable transportation infrastructure. In maintenance operations, obtaining limits of homogeneous sections is a key problem because appropriate segmentation can help yield a more cost effective M&R plan. The purpose of this study is to present the result of investigation on various research works and to suggest the direction of developing an enhanced methodological framework. Existing approaches for pavement segmentation was explored through a literature review and data analysis. Autocorrelation tests, change-point approaches, a Bayesian method, and a hidden Markov model were performed using pavement condition data. Future work directions were suggested to develop a segmentation method capable of handling the issues found in the study.

Table of Contents

List of Tables	ix
List of Figures	x
Chapter 1: Introduction	1
Chapter 2: Literature Review: Existing Approaches for Pavement Segmentation	6
2.1 Cumulative Difference Approach	6
2.2 Bayesian Approach	9
2.3 Fuzzy c-mean clustering	11
2.4 Wavelet Transform	13
2.5 CART	13
2.6 MINSSE	14
2.7 Discussion	15
Chapter 3: Methodology and Data Analysis	17
3.1 Test Data for Analysis	17
3.2 Autocorrelation Test of Road Performance Data	19
3.2.1 Methods for Testing Autocorrelation	19
3.2.1.1 ACF and PACF	19
3.2.1.2 Durbin-Watson Test	20
3.2.2 Data Analysis for testing autocorrelation	21
3.2.2.1 Condition Score	21
3.2.2.2 Ride Score	24
3.2.3 Discussion for Testing Autocorrelation	27
3.3 Off-the-shelf Codes	28
3.3.1 Change-point Detection Packages in R	29
3.3.1.1 R package: changepoint	29
3.3.1.2 R package: bcp	30
3.3.1.3 R package: ecp	31
3.3.1.4 Discussion of Change-point Detection Packages in R	33

3.3.2 Hidden Markov Model in MATLAB.....	34
3.3.2.1 Hidden Markov Model.....	34
3.3.2.2 A toolkit for Hidden Markov Model in MATLAB.....	35
3.3.2.3 Discussion of Hidden Markov Model in MATLAB.....	37
3.4 Implement A Bayes Approach.....	38
3.4.1 Data Used in Analysis.....	38
3.4.2 Data Analysis.....	39
3.4.2.1 Bayesian Approach with Unknown Common Variance.....	39
3.4.2.2 Bayesian Approach with First-order Autoregressive Process (AR(1)).....	43
3.4.3 Discussion.....	44
Chapter 4: Conclusion and Future Study.....	46
References.....	49

List of Tables

Table 1:	Guidelines for interpreting Bayes factors (After Jeffreys, 1998)	11
Table 2:	Descriptive Information of Test Data	18

List of Figures

Figure 1:	Consequences of segmentation, adapted from Cafiso and Graziano, 2012 (after Acurio, 2014)	3
Figure 2:	Reasoning of segmentation (After Yang et. al., 2009).	4
Figure 3:	Concept of cumulative difference approach to analysis unit delineation (after AASHTO, 1993)	8
Figure 4:	An example of optimal solutions for FCM algorithm (After Yang et. al., 2009)	12
Figure 5:	Example of the CART result: (a) original tree; (b) sub-tree (After Misra et. al., 2003)	14
Figure 6:	Comparison of different segmentation method: CDA, Bayesian, and MINSSE (After Cafiso et. al. 2012).....	15
Figure 7:	SH0046 on TxDOT's Statewide Planning Map (Blue Line).....	17
Figure 8:	Test Data Plot (Left) and Histogram of Condition Score (Right).....	19
Figure 9:	Plots of ACF (a) and PACF (b) with respect to condition score for whole segment	22
Figure 10:	Segmentation result to subset X1 and X2	23
Figure 11:	Plots of ACF ((a), (c)) and PACF ((b), (d)) with respect to condition score for partial segments	24
Figure 12:	Ride score for the test sections of SH0046	25
Figure 13:	Plots of ACF (a) and PACF (b) with respect to ride score for whole segment	26
Figure 14:	Plots of ACF ((a), (c)) and PACF ((c), (d)) with respect to ride score for partial sections	27

Figure 15:	Resulting Plots of the changepoint Package in R.	29
Figure 16:	Resulting Plots of the bcp Package in R.	31
Figure 17:	Resulting Plots of the ecp Package in R.	32
Figure 18:	Graphical representation of HMM.....	34
Figure 19:	HMM results with varying the number of states equals to (a) 3; (b) 4; (c) 5; (d) 8.....	36
Figure 20:	(a) Plot and (b) histogram of the ride score data.....	39
Figure 21:	Result plot of the first run and posterior probabilities for a change point	40
Figure 22:	Result of detecting multiple change points using the Bayesian approach	41
Figure 23:	Result of the Bayesian algorithm illustrating a multi modal case using the 7 th segment data.....	42
Figure 24:	Bayes factor for ‘change’ vs. ‘no change’ vs. various values of φ .	44
Figure 25:	Representation of 5-fold cross validation	47
Figure 26:	5-fold cross validation for the Bayesian approach.....	48

Chapter 1: Introduction

Pavement maintenance and rehabilitation (M&R) is important for transportation agencies to have a sustainable transportation infrastructure. Pavement maintenance consists of routine and preventive activities such as filling cracks, patching, chip seal, and so on. Pavement rehabilitation includes actions such as overlay and partial to complete reconstruction that increase the structural capacity of pavement. Due to the size of road networks, M&R is one of the major investments in a transportation system. Accordingly, planning M&R is a problem that challenges decision makers because they need to determine which pavement road section has to be treated, when and how that treatment should be conducted. In addition, the decision making process must take into account budget limitations, meet specific goals for maintaining pavement performance, and allocating budgets to maximize cost effectiveness (Hass et. al., 1994)

In the Texas Department of Transportation (TxDOT), the Pavement Management Information System (PMIS) has been operated since the early 1990s to support pavement related decision making processes by storing, retrieving, analyzing and reporting information (TxDOT, 2003). Currently, the information managed in 0.5 mile data collection section, thus, it can be used in analysis such as condition estimation and maintenance needs estimation for administrative level. However, for district level project selection, the half mile section data are restrictive because typically projects are of any length that combines multiple half mile sections. Therefore, instead of using the half mile data collection section, using the management section consisting of homogeneous sections is necessary. For that reason, obtaining the limits of homogeneous sections becomes a key problem in pavement management.

Scullion and Smith (1997) presented three options for selecting the limits. Firstly, use the control sections which were designed and constructed under identical conditions. Second is to use the limits proposed by pavement engineers. Lastly, the cumulative difference approach (CDA) using pavement performance indices such as ride and condition scores could be used to delineate sections (Scullion and Smith, 1997). Each approach has its limitations. For example, the control sections are not guaranteed to be homogeneous after maintenance and rehabilitation works are done because typically M&R projects cover partial sections within the control sections. Also, defining limits using engineering judgement seems too demanding to cover large road networks in the case of the second option. Moreover, researchers have agreed that CDA has limitations which will be examined in later chapter in this report. On top of that, there has been no official method for obtaining a homogenous management section. Therefore, it would seem that further investigation is needed in the area of finding homogeneous sections, namely, segmentation.

Segmentation methods can be split into three categories: fixed length segmentation, dynamic segmentation, and static segmentation. In case of the fixed length segmentation, the segment limits originate from fixed features and are kept constant over time. For instance, the aforementioned approach to determine the limits using control sections is an example of fixed length segmentation. In the case of dynamic segmentation, a decision for segment's boundaries is based on the homogeneity of pavement sections' attributes including ride quality and conditions. In this setting, the boundaries vary as the attributes changes over time. Static segmentation is similar to dynamic segmentation, except the limits of segments are kept for a certain time to make management easier (Bennett, 2004).

Appropriate segmentation can help yield a more cost effective M&R plan. Figure 1 illustrates the advantage of having proper segmentation. If one defines a whole section as a single segment without taking account into a change in pavement attribute, treatment must be applied on pavement sections with good condition as shown in the middle of Figure 1. Conversely, if one properly identifies two segments as in the bottom of Figure 1, the segment requiring maintenance will be treated, which leads to cost effective M&R plans (Cafiso and Graziano, 2012).

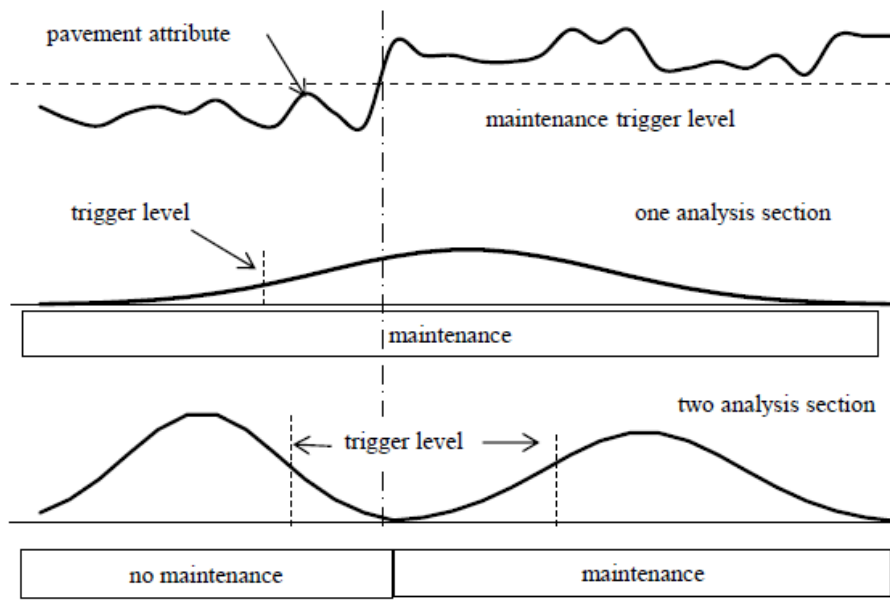


Figure 1: Consequences of segmentation, adapted from Cafiso and Graziano, 2012 (after Acurio, 2014)

Another example is demonstrated in Figure 2. The historical project limits shown in (a) in the figure might not yield homogeneous pavement conditions because the pavement sections might experience different traffic and maintenance history. Thus, if one keeps the historical project limits as (b) in the figure, the resulting segmentation does not take into account the changes in pavement conditions. Under a better segmentation

scheme, two segments are able to become three segments whose conditions are more homogeneous within each segment as shown in (c). As a result, more cost effective M&R plan can be established (Yang et. al., 2009).

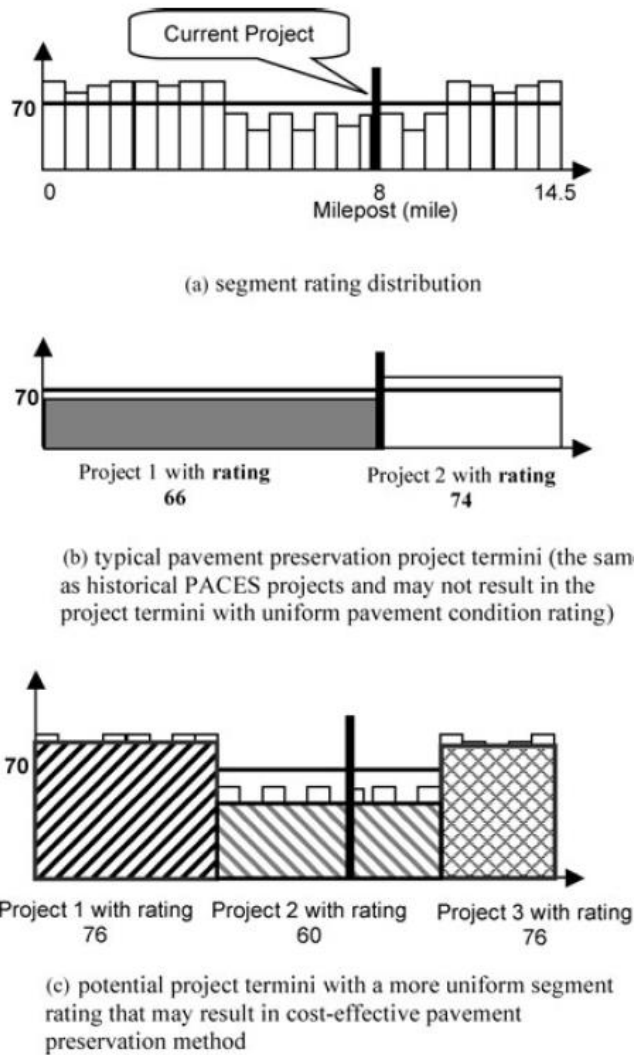


Figure 2: Reasoning of segmentation (After Yang et. al., 2009).

The purpose of this study is to present the result of investigation on various research works and to suggest the direction of developing an enhanced methodological framework. The remainder of this report is organized as follows. Chapter 2 presents the

result of literature review regarding research studies on the segmentation methods of pavement sections. In Chapter 3, tests of autocorrelation on pavement condition indices were conducted. In addition to that, I explored off-the-shelf tools available for detecting a change point in R and MATLAB, and also implemented a Bayesian approach. I conclude this report with suggesting future work direction in Chapter 4.

Chapter 2: Literature Review: Existing Approaches for Pavement Segmentation

Throughout the literature review on segmentation methods of highway application, it was found that there have been various research works to develop an approach to identify homogeneous segments.

2.1 CUMULATIVE DIFFERENCE APPROACH

The Cumulative Difference Approach (CDA) is by far the most popular method for segmentation in the pavement management sector. One of the reasons for the popularity is that the method is included in the AASHTO pavement design guide (AASHTO, 1993) used in worldwide as a pavement design guideline. Another reason is that the method is straightforward and powerful as stated in the AASHTO guide (AASHTO, 1993).

Figure 3 shows the overall concept of CDA. As shown in Figure 3 (a), there are three unique constant values r_1 , r_2 and r_3 with three intervals 0 to x_1 , x_1 to x_2 , and x_2 to x_3 , respectively. The cumulative area at x can be calculated as the following integral:

$$A = \int_0^{x_1} r_1 dx + \int_{x_1}^x r_2 dx$$

The cumulative area of the average project response can be calculated by following equations:

$$A_x = \int_0^x r dx$$

with

$$\bar{r} = \frac{\int_0^{x_1} r_1 dx + \int_{x_1}^{x_2} r_2 dx + \int_{x_2}^{x_3} r_3 dx}{L_p} = \frac{A_T}{L_p}$$

and therefore

$$\overline{A}_x = L_P \times A_T$$

The cumulative difference variable Z_x is determined as the following relationship. In Figure 3 (b), Z_x is illustrated as the difference between cumulative areas at x .

$$Z_x = A_x - \overline{A}_x$$

When Z_x is plotted over the length of project as illustrated in Figure 3 (c), the boundary location can be determined by the location where the slope of Z_x function changes, for example from negative to positive or vice versa (AASHTO, 1993).

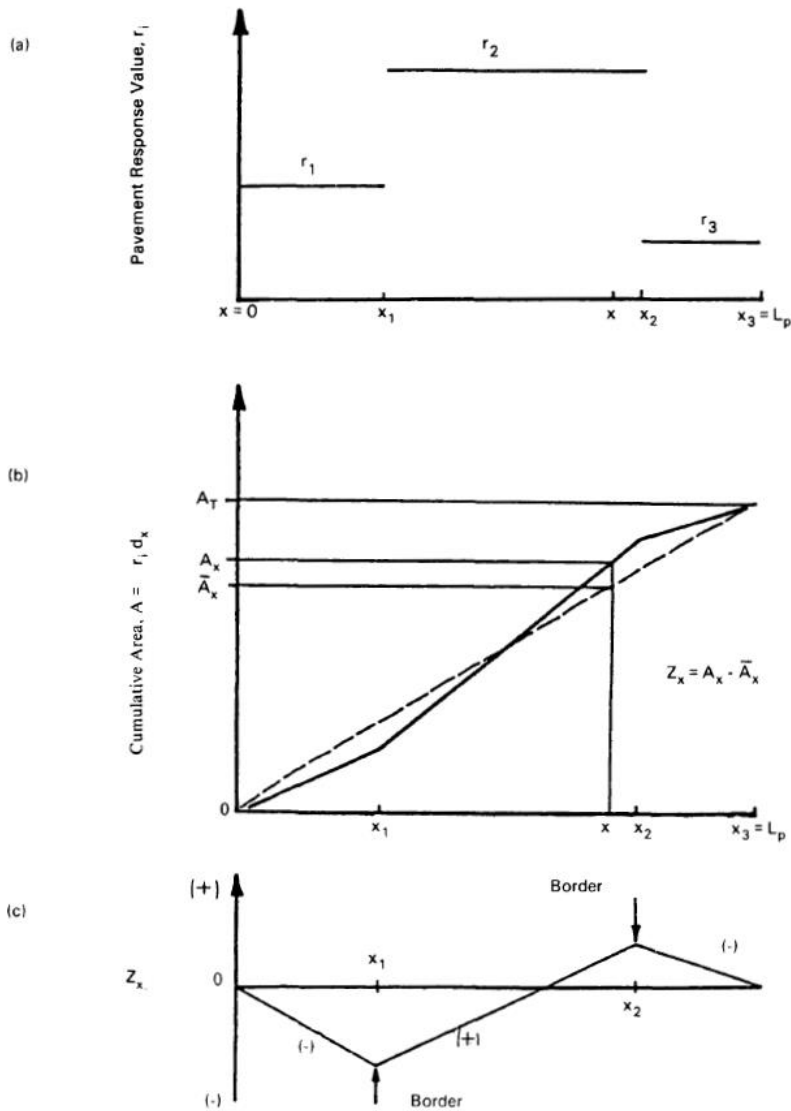


Figure 3: Concept of cumulative difference approach to analysis unit delineation (after AASHTO, 1993)

Various research works discussed limitations of the CDA and suggested modified procedures of CDA or new methods. Misra and Das (2003) discussed limitations of CDA. In the case of more than one homogeneous sections with different mean levels consecutively exist above or below the mean horizontal line, the CDA fails to delineate

those section because the sign of Z_x does not change. Also, they mentioned CDA has no control over the number of homogenous section, and minimum section length is not chosen by user. They suggested a CART algorithm as an improved method.

Divinsky et. al. (1997) recommended a modification of the CDA procedure by taking into account statistically homogeneous scatter characteristics such as standard deviation, range, etc. to overcome the limitation of significant sensitive to the existing change in the means of segments.

Ping et. al. (1999) introduced a procedure for automated segmentation of pavement rut data using CDA. They developed a multipass SAS program that runs on the entire data set in the first set and then iterates process with resulting segments to obtain subsegments until the program produce the same segments as the previous pass. In the program, two user specified constraints such as a minimum segment length and a minimum difference in mean are incorporated. Those constraints were compared by the sum of squared errors (Ping et. al., 1999). Kennedy, Shalaby, and Cauwenberghe (2000) conducted the CDA on IRI data with a similar procedure as the study of Ping et. al. (1999). Cafiso, Di Graziano (2012) also used the same procedure to compare their method MINSSE and the CDA.

Thomas (2005) argued that the CDA is mostly a graphical method to detect the homogeneous sections, and it is not suitable for narrowly spaced measurements. Also, the CDA always suggests at least two segments unless all measurements in a given series are identical. He introduced a Bayesian approach which will be discussed in the next section.

2.2 BAYESIAN APPROACH

Thomas (2003) presented a method to detect a change in the mean, in the variance and/or in the autocorrelation of a series using a Bayesian approach that allows

communicating the existence and possible location of a change point in terms of probabilities. The author emphasized that the method requires no prior knowledge and distributional assumption. Later, Thomas (2005) introduced Box-Cox transformations to meet the normality assumption of observations, and a heuristic algorithm to detect multiple change points to overcome the limitation of at most one change point algorithm. These two studies based on his dissertation thesis (Thomas, 2001). Detailed statistical proofs are presented in the thesis, so here I introduce the basic concept of a Bayesian approach.

A general approach of a Bayesian change point is introduced in Thomas' thesis as follows. A sequence of random variables, x_1, \dots, x_n , is divided into subsequences $x_1, \dots, x_r; x_{r+1}, \dots, x_n$ by a change point r , where $1 \leq r < n$. M_0 indicates the model with no change in underlying parameters and its joint density can be expressed as $p(x_1, \dots, x_n | M_0)$. Meanwhile, a model with change in one or more of the parameters at r is denoted as M_r and its joint density is $p(x_1, \dots, x_n | M_r)$. Using Bayes theorem, the posterior probability of a model is,

$$p(M_r | x_1, \dots, x_n) = \frac{p(x_1, \dots, x_n | M_r)p(M_r)}{\sum_{all\ r'} p(x_1, \dots, x_n | M_{r'})p(M_{r'})}$$

$$\propto p(x_1, \dots, x_n | M_r)p(M_r)$$

Then two model comparisons are interested. One is comparing the models M_1, \dots, M_{n-1} to specify change in parameters occurs at $r = 1, \dots, n - 1$. Another is comparing some or all models have a change at r with M_0 to test whether a change occurs at all. Comparing two models, where change point at r and s , respectively, can be conducted by following Bayes factor. In general, Bayes factors provide a way of quantifying the evidence based on data in favor of a null hypothesis.

$$\frac{\frac{p(M_r|x_1, \dots, x_n)}{p(M_s|x_1, \dots, x_n)}}{\frac{p(M_r)}{p(M_s)}} = \frac{p(x_1, \dots, x_n|M_r)}{p(x_1, \dots, x_n|M_s)} = B_{rs}$$

Comparing the hypothesis of no change vs. a change in series can be done using,

$$\frac{1 - p(M_0|x_1, \dots, x_n)}{p(M_0|x_1, \dots, x_n)} / \frac{1 - p(M_0)}{p(M_0)} = \sum_{r=1}^{n-1} B_{r0} \frac{p(M_r)}{1 - p(M_0)}$$

A Bayes factor under the assumption that the numerator and the denominator are identical can be interpreted using the guideline given in Table 1. (Thomas, 2001)

Bayes factor	Interpretation
> 100	Decisive evidence for H _A
30–100	Very strong evidence for H _A
10–30	Strong evidence for H _A
3–10	Substantial evidence for H _A
1–3	Anecdotal evidence for H _A
1	No evidence

Table 1: Guidelines for interpreting Bayes factors (After Jeffreys, 1998)

2.3 FUZZY C-MEAN CLUSTERING

Yang, Tsai and Wang (2009) developed a spatial clustering algorithm using Fuzzy C-mean Clustering (FCM). The algorithm minimizes the pavement condition rating variation in each project while take into account minimum length, costs, barrier, and pavement surface type of a project.

In order to accomplish the goal, the algorithm uses two objective functions. One is for minimizing rating variation, and another is for minimizing costs for projects. Together with constraints, the optimal result can be achieved. The optimization process is

repeated in the range of cluster numbers. Among the multiple results of optimal number of clusters, the best segmentation would be selected based on the cost objective function.

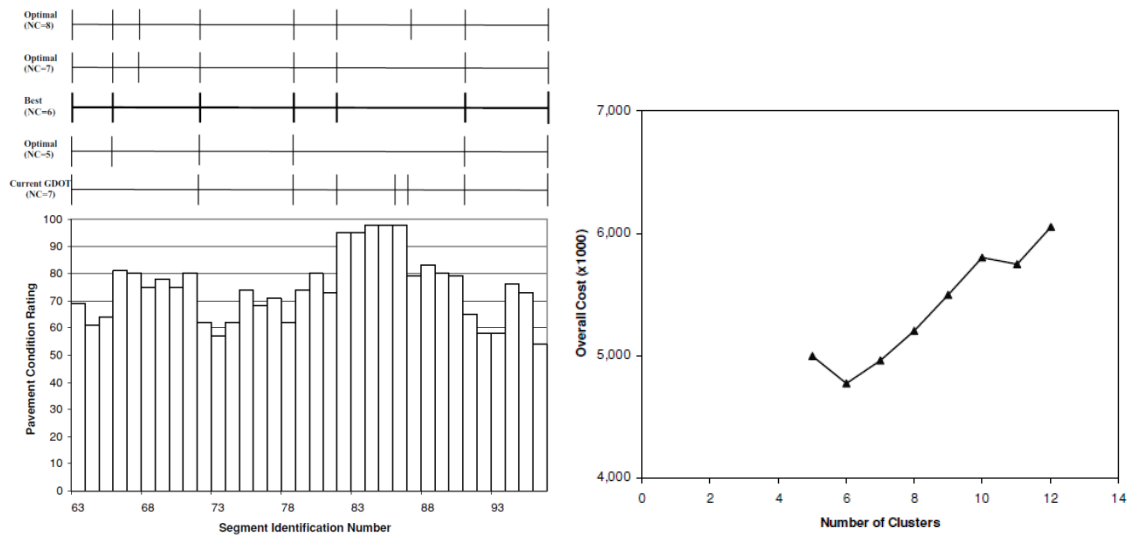


Figure 4: An example of optimal solutions for FCM algorithm (After Yang et. al., 2009)

The left hand side of Figure 4 shows partitions found when applying the FCM procedure to SR 10 in Georgia, with 5, 6, 7 and 8 clusters. The right hand side shows the associated cost for each segmentation; from this we see that the best segmentation case is when the number of cluster is equal to 6 based on the minimizing cost criteria with satisfying all constraints.

This method offers a way to cluster sections with optimization scheme that includes costs. Thus, conceptually, it provides a better solution than a simple CDA; however, the method does not provide any statistical inference.

2.4 WAVELET TRANSFORM

An algorithm based on wavelet transforms for automated segmentation was presented by Cuhadar et. al. (2002). The properties of wavelet transform such as de-noising and singularity detection were used to delineate sections with respect to the pavement condition data. The original data is transformed to a smoother waveform by using de-noising, and then, singularity detection was applied on the smoothed data. The algorithm results in the pavement condition data into regions which have similar characteristics (Cuhadar et. al., 2002).

Boroujerdian et. al. (2014) also used wavelet theorem for the dynamic segmentation. In the study, based on the wavelet theory, the length of high crash road segments is identified by converting accident data to the road response signal.

Wavelet transformation seems outperform the CDA because this approach overcomes the sensitivity to small variability in data by de-noising. However, the method of singularity detection seems able to identify only sudden changes in level of data.

2.5 CART

Misra et. al. (2003) proposed a method using classification and regression trees (CART) (Breiman et. al., 1984). The algorithm produces a binary tree through the exhaustive search to find the point that minimizes the sum of squared error. By recursive binary splitting, the original tree will be produced as shown in Figure 5 (a). Once the original tree is produced, based on constraints such as a minimum section length and a number of sections, the tree is reduced by merging dividend sections in the original tree. In Figure 5 (b), the resulting sub-tree is illustrated which indicates 8 delineated sections.

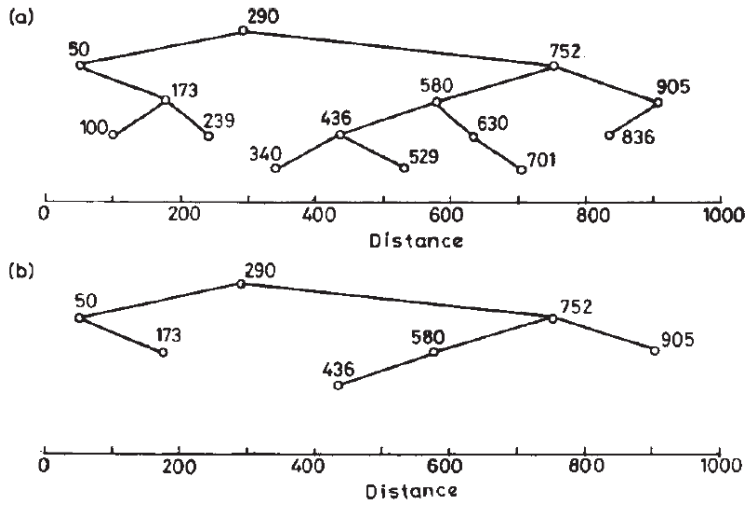


Figure 5: Example of the CART result: (a) original tree; (b) sub-tree (After Misra et al., 2003)

The proposed algorithm provides a simple and fast solution for segmentation without any assumption of the distribution of data. The limitations of CDA are overcome by constraining the minimum length of section and choosing the number of segments. Nonetheless, this approach does not produce the optimal solution because it uses recursive binary trees as approximations.

2.6 MINSSE

Cafiso et al. (2012) introduced the minimum sum of squared error (MINSSE) method. The method is to find the minimal sum of squared error (SSE) of partitions as following:

$$SSE_k = \sum_j^{k+1} \sum_{i \in S_j} (x_i - \bar{x}_{S_j})^2$$

where, k is the number of segments. S_j a set of element in j^{th} segment.

Once change points minimizing the SSE determined under a given minimum length of segment, t-tests are conducted to check if adjacent segments meet the criteria of a minimum difference and those segments are combined if the test fails. The authors compared MINSSE with the CDA and the Bayesian approach, and concluded that their method resulted in similar segmentation to the Bayesian approach although their method is less complex to conduct. Figure 6 is the comparison graph of three approaches such as CDA, Bayesian, and MINSSE using rut data.

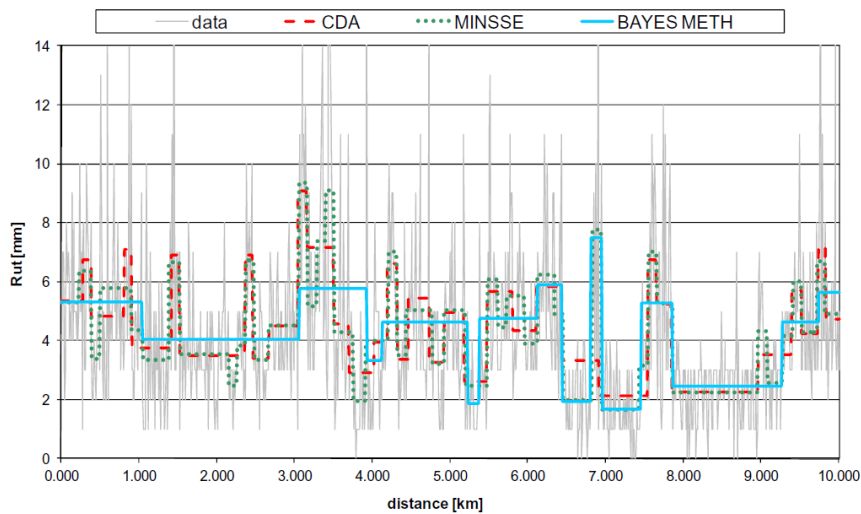


Figure 6: Comparison of different segmentation method: CDA, Bayesian, and MINSSE (After Cafiso et. al. 2012)

2.7 DISCUSSION

The majority of existing studies segment road data based on pavement performance data such as IRI, skid, rut, etc. There have been several safety related research works; however, only one work based on wavelet theorem was introduced because an analysis on crash count data has different characteristics-a segment determined by criteria that a fixed length with a certain crash counts.

There seems to be general agreement on the limitations of the CDA. Therefore, some studies modified the CDA to improve it and the other studies suggested new methods that showed better performance. The developed methods commonly adopted constraints, such as a minimum section length, a minimum difference, a number of segments, etc. to overcome the limitations of the original CDA.

Most studies have focused on delineating segment based on different mean levels of segments. Few studies have attempted to develop a method which takes into account variance and autocorrelation. In addition, no studies have explored on multivariate data. For example, no method can conduct segmentation based on rut and skid data simultaneously. Thus, it would be of interest to develop a method that can implement aforementioned gaps.

Chapter 3: Methodology and Data Analysis

In this chapter, tests of autocorrelation on pavement condition indices were conducted. And then, off-the-shelf tools for detecting a change point in R and MATLAB were explored. Lastly, a Bayesian approach implemented using R.

3.1 TEST DATA FOR ANALYSIS

In Texas, the road network is huge and variety. Thus, it is not reasonable at all to select one road section for representing whole network. However, for testing various methods of change point problem, I picked a road section whose length is relatively short, but not too short section. Also, the picked section should have not too simple but complex enough profile of a performance measure.

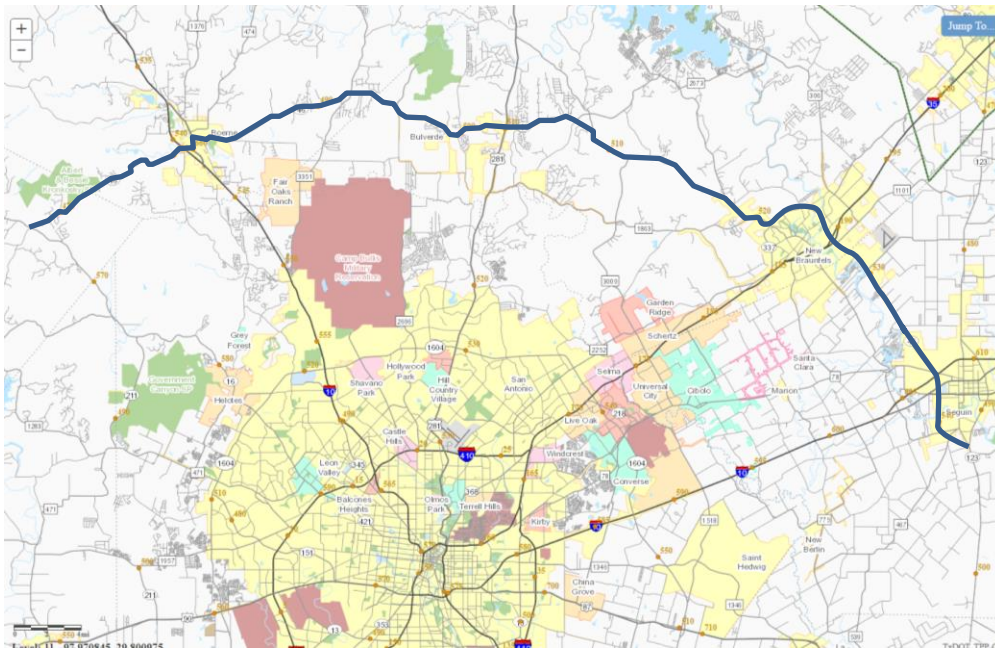


Figure 7: SH0046 on TxDOT's Statewide Planning Map (Blue Line)

Candidate sections were screened within TxDOT's San Antonio district. After inspecting map and data, State Highway 46 (SH0046K) was selected as the test data.

Figure 7 shows the map of SH0046 in San Antonio district, Texas. It is state highway links northern outskirts of San Antonio from west to east by 2-4 through lanes. Some details about the section are described in Table 2, and the histogram of condition scores for the test data is given on the right hand side of Figure 8. When it comes to the average condition score, it is well maintained road section. However, as shown in Figure 8 on the left, whole section cannot be treated as a homogenous section because some sections show relatively lower condition scores. There were zero values in condition scores which indicate 3 missing data points. I left these missing data on purpose to see how the testing methods work on circumstance of having missing data in series, and study further regarding missing data treatment.

Year of Data	2013
Location	TRM 468 + 0 – TRM 542 + 0.7
Length	67.4 miles
Number of Sections	139
Pavement Type	Asphalt Concrete
Condition Score Average	89.63
Condition Score Max.	100.00
Condition Score Min.	0.00
Condition Score Standard Deviation	19.14

Table 2: Descriptive Information of Test Data

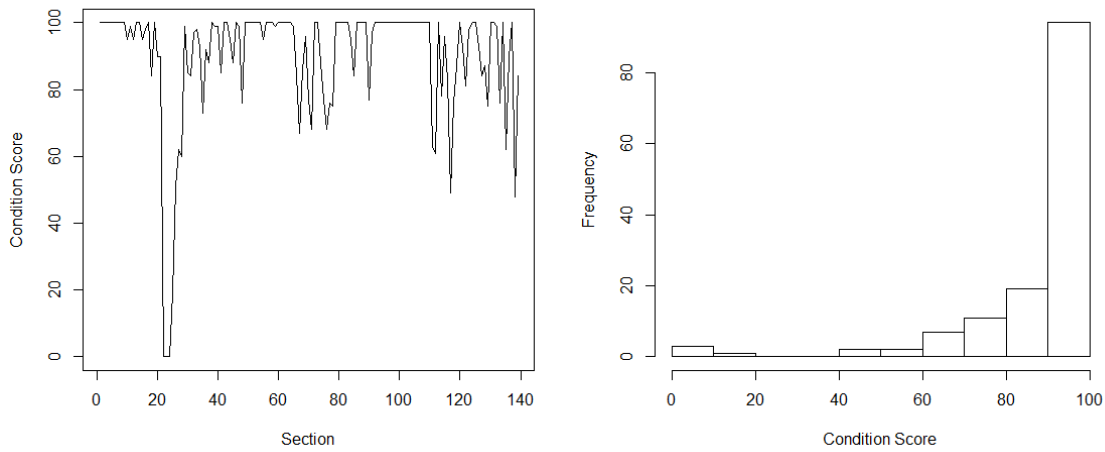


Figure 8: Test Data Plot (Left) and Histogram of Condition Score (Right)

3.2 AUTOCORRELATION TEST OF ROAD PERFORMANCE DATA

The road performance data are likely to have autocorrelation in nature because adjacent sections experience similar traffic and environmental conditions. Therefore, the assumption of independence might not be reasonable even though it derives simpler model. The test data including condition score and ride score were examined to find whether or not autocorrelation needs to be considered.

3.2.1 Methods for Testing Autocorrelation

There are various methods for testing whether data exhibits autocorrelation. In this study, visual inspections of the Autocorrelation Function (ACF) plot and Partial Autocorrelation Function (PACF) plot to identify the existence of autocorrelation and the Durbin-Watson test for test AR(1) process were used.

3.2.1.1 ACF and PACF

The Autocorrelation Function (ACF) is to describe the properties of a stationary stochastic process, and the theoretical ACF is defined as follows:

$$\rho(k) = \text{corr}(x_t, x_{t-k}) = \frac{\text{cov}(x_t, x_{t-k})}{\sqrt{\text{var}(x_t)}\sqrt{\text{var}(x_{t-k})}}$$

The sample autocorrelation, the estimator for the ACF, can be derived from the sample autocovariance:

$$\hat{\rho}(k) = \frac{\sum_{t=1}^{n-k} (x_t - \bar{x})(x_{t-k} - \bar{x})}{\sum_{t=1}^n (x_t - \bar{x})^2}$$

A plot of the sample autocorrelation is called a correlogram that provides a way to judge which stochastic process would be a suitable model for data (Durbin and Watson, 1951).

The ACF incorporates the intermediate linear correlations between x_t and x_{t-k} , for example, $\text{corr}(x_t, x_{t-2})$ takes into account for the correlations between x_t , x_{t-1} and x_{t-1} , x_{t-2} . The partial autocorrelation function (PACF) is another function that provides additional information about autocorrelation to represent the net correlation between x_t and x_{t-k} by eliminating the intermediate linear relationship. The Durbin-Levinson algorithm can be used to estimate PACF (Levinson, 1947; Durbin, 1960). In practice, by looking at the estimated PACF plot, it is possible to make a decision if an autoregression may be a suitable for data. For example, in case of the AR(p) model, estimated partial autocorrelation coefficients should be significant at the p-th lag. (Prado and West, 2010)

3.2.1.2 Durbin-Watson Test

The Durbin-Watson test (Durbin and Watson, 1951) is conducted to test whether the residuals from a linear regression are independent. A general assumption of a linear regression is the independence of each error term. If there is a relationship between

neighboring error terms, for example error at t-1 and error at t, it is called first order autocorrelation AR(1). The following relationship is the AR(1) process:

$$\varepsilon_t = \rho\varepsilon_{t-1} + \mu_t$$

where, ρ is autocorrelation parameter, $-1 \leq \rho \leq 1$.

The hypotheses test is

$$H_0: \rho = 0$$

$$H_a: \rho > 0$$

The test statistic is

$$D = \frac{\sum(e_t - e_{t-1})^2}{\sum e^2}$$

The test statistic D ranges from 0 to 4. When there is no first order autocorrelation, we expect D to be close to 2. When D is smaller than 2, positive autocorrelation is expected. D greater than 2 indicates negative autocorrelation. By using the provided table, one-sided test for positive autocorrelation can be decided as follows:

If $D < D_L$ then reject H_0

If $D > D_U$ then do not reject H_0

If $D_L < D < D_U$ then the test is inconclusive

More detailed explanation and the table can be found in references (Kutner, 2005; Montgomery et. al., 2013)

In this study, the Durbin-Watson test implemented in R package lmtest was used to analyze the road performance data.

3.2.2 Data Analysis for testing autocorrelation

3.2.2.1 Condition Score

When full sections of sample data were used in analysis for testing autocorrelation, ACF showed exponential decay and PACF showed a significant lag at 1

(see Figure 9). Those are the evidences of AR(1) process. In addition, Durbin-Watson test result showed that the p-value is $1.678e-15$, which means that the null hypothesis is rejected at 95% significant level. All in all, it can be concluded the full sections follows AR(1) process.

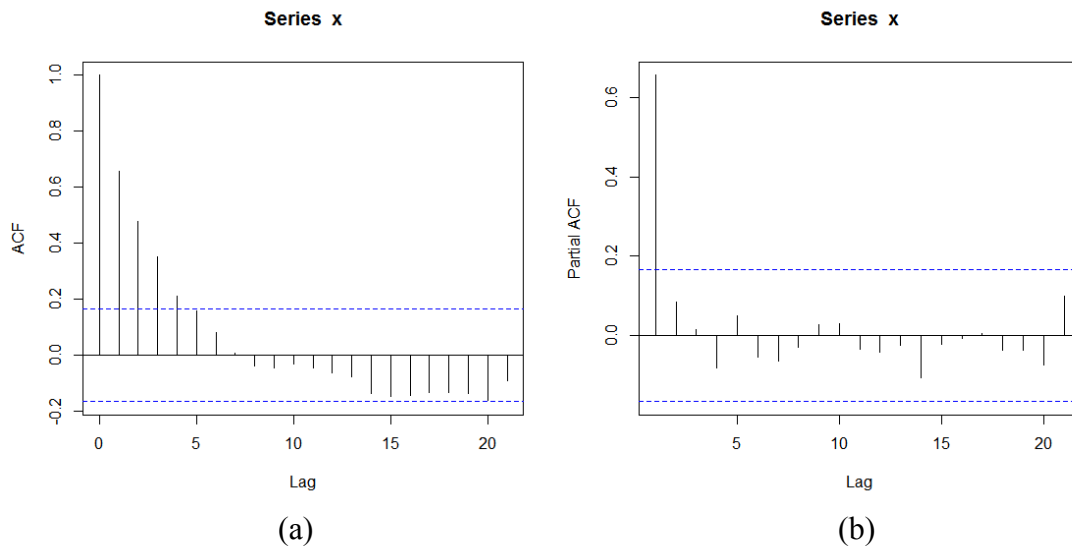


Figure 9: Plots of ACF (a) and PACF (b) with respect to condition score for whole segment

However, it seems likely that the high autocorrelation is due to many sections of road having continuous stretches where the condition score takes the maximum value 100 (see Figure 8). By using the segmentation result of R package ‘changeoint’ as can be seen Figure 10, two non-maximum parts X1 and X2 were taken out and tested separately to verify the effect of these continuous maximum values.

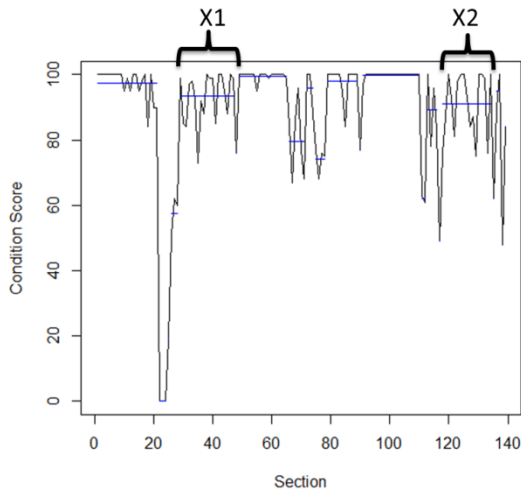


Figure 10: Segmentation result to subset X1 and X2

I found that, both subset X1 and X2 did not show autocorrelation. As shown in Figure 11, the plots of ACF and PACF demonstrated no evidence of autocorrelation. They had rather properties of random walk. Also, the result of Durbin-Watson test verified this fact because the p-value for X1 and X2 were 0.4747 and 0.1971, respectively. That is, the null hypothesis cannot be rejected due to the large p-value, which means the autocorrelation parameter ρ is equal to zero.

The resulting non-autocorrelation behavior of the subsets might be caused by the nature of the test data. As a different dataset may lead to a different result, I took another segment from Interstate Highway 35 (IH35) in San Antonio district, and then the same test procedure was conducted as previous. The results of test using IH35 data were similar to the previous test. In the case of test using full sections of data showed autocorrelation due to the continuous maximum values. Among three subsets tested for sections contain non-maximum values, 2 subsets gave no evidence of autocorrelation, but 1 subset slightly showed evidence of autocorrelation. Figures and values for the results of

ACF, PACF, and Durbin-Watson test were omitted because they are very similar to the previous test results.

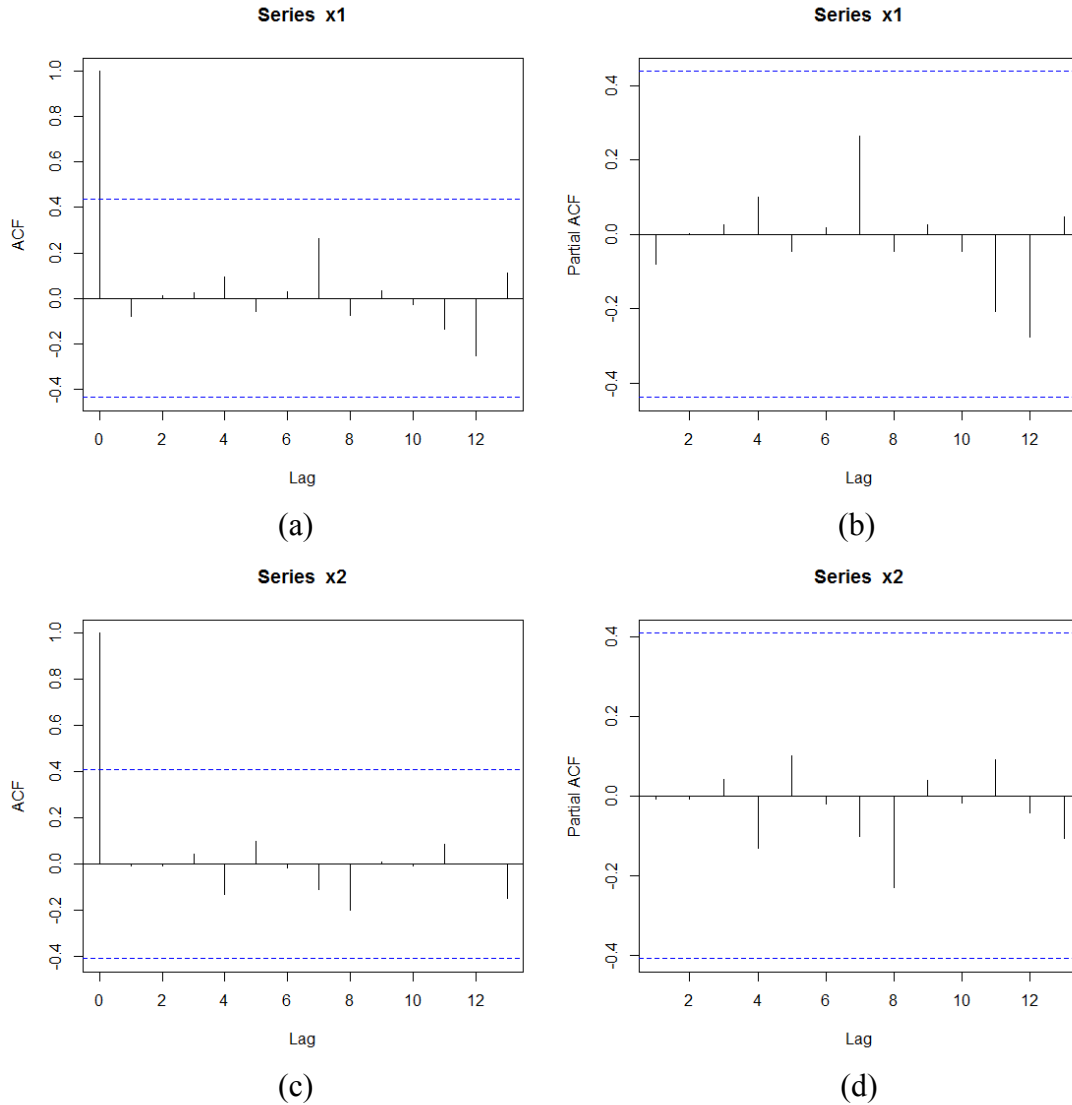


Figure 11: Plots of ACF ((a), (c)) and PACF ((b), (d)) with respect to condition score for partial segments

3.2.2.2 Ride Score

Condition score data showed strong evidence of autocorrelation in the case of using whole segment of the sample data based on ACF, PACF and Durbin-Watson test.

However, the result might be caused by the fact that many sections in condition score data had the maximum value (100) in a row. When partial sections without the maximum condition score were used in the same analysis, the results indicated no autocorrelation. Therefore, it is not good idea to have condition score as a sample data to test autocorrelation. Ride score was considered as a substitute measure for the analysis. Ride score is another index in PMIS to indicate the roughness of pavement surface, ranged in values from 0 to 5. Figure 12 shows ride scores for the same segment of SH0046. As opposed to the condition score data, it does not show any consecutive trend of the maximum value.

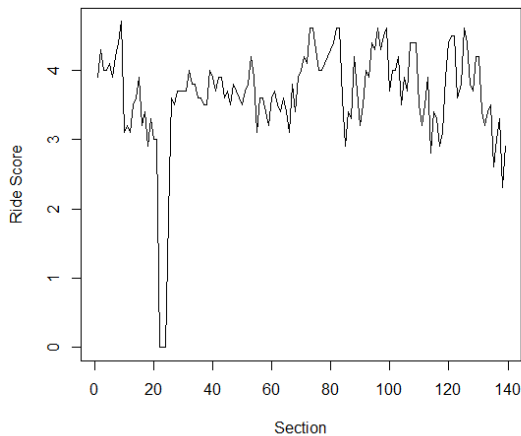


Figure 12: Ride score for the test sections of SH0046

The analysis was done by the same procedure. First, the whole segment was used for the test. As shown in Figure 13, ACF showed exponential decay and PACF showed a significant lag at 1. In addition, Durbin-Watson test result showed that the p-value is $2.2e-16$, which means that the null hypothesis is rejected at 95% significant level. All in all, it can be concluded the full sections follows AR(1) process.

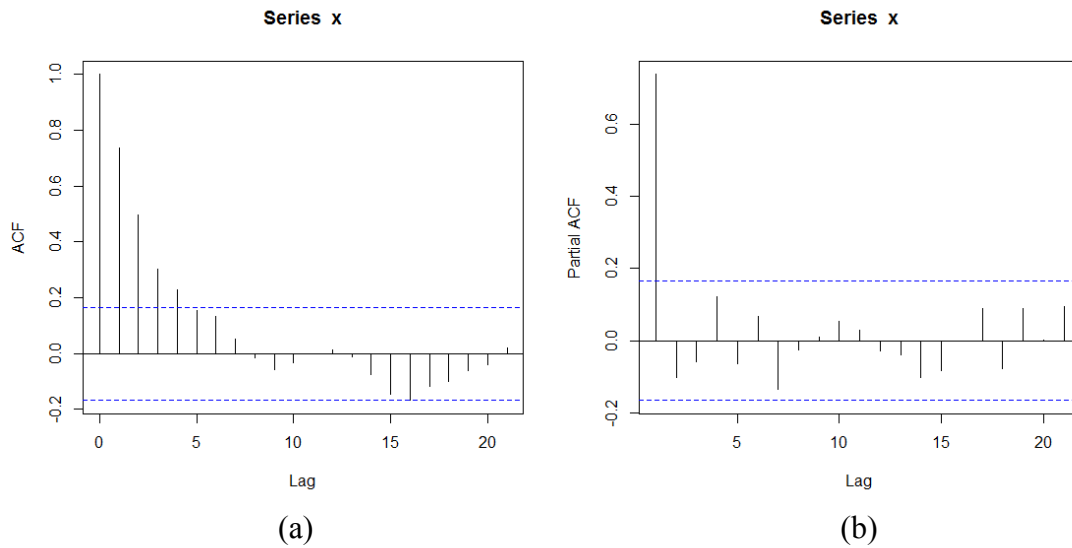


Figure 13: Plots of ACF (a) and PACF (b) with respect to ride score for whole segment

As a next step, by using the ‘change point’ results, two subsets of data X1 and X2 were tested as previously. As for the subset X1, ACF showed no exponential decay, and PACF had no significant lag. In this end, while visually, the subset X1 seemed not to have significant autocorrelation as it is evident from Figure 14. The Durbin-Watson test result indicated that the series shows autocorrelation. The test statistic D was 1.1607 and the p-value was 0.00104. For the subset X2, ACF plot showed moderate exponential decay and PACF had significant lag at 1, which are the properties of AR(1). The Durbin-Watson test also provided an evidence of AR(1) as the test statistic D was 1.2642 and the p-value was 0.01495.

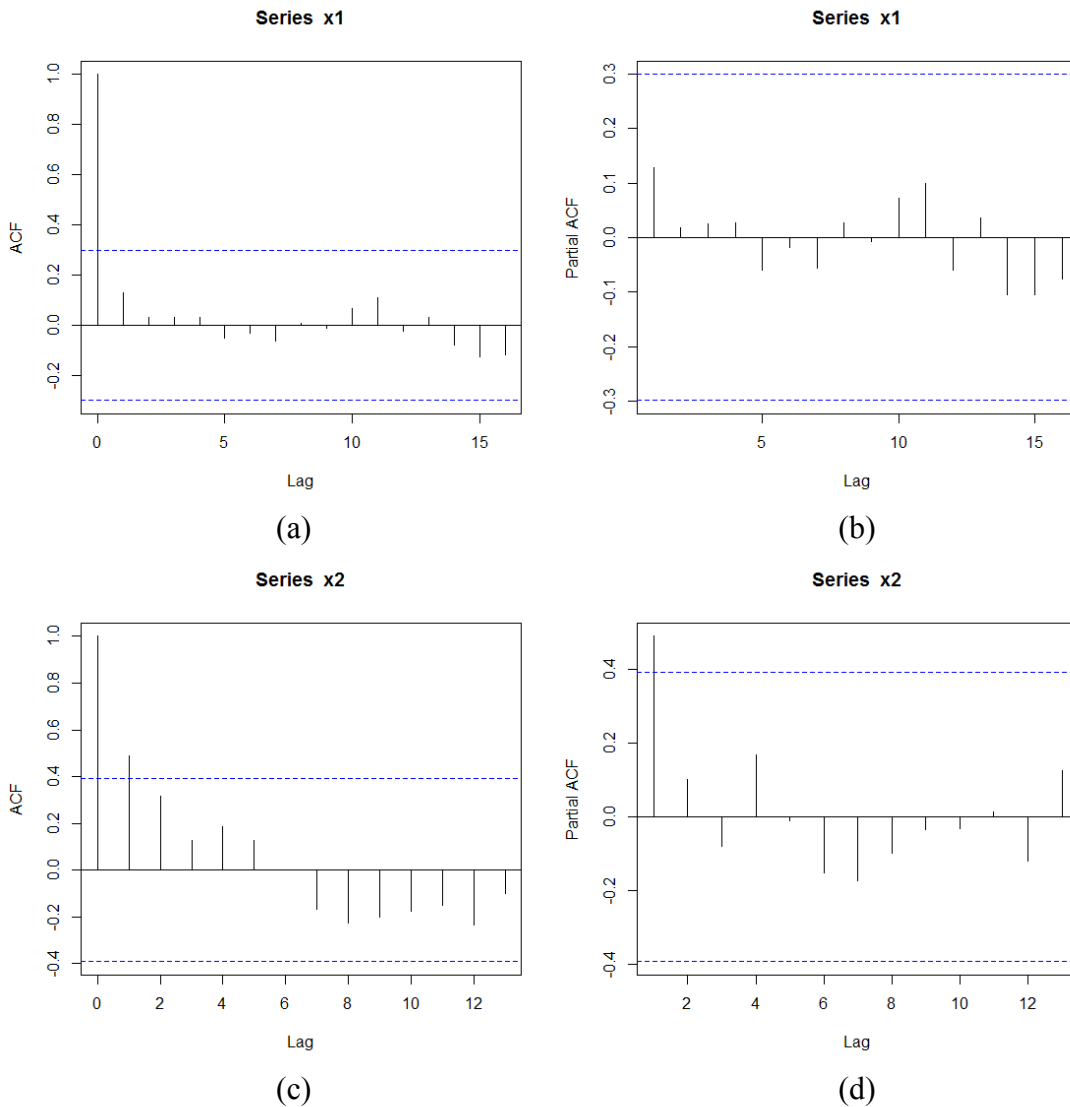


Figure 14: Plots of ACF ((a), (c)) and PACF ((c), (d)) with respect to ride score for partial sections

3.2.3 Discussion for Testing Autocorrelation

Condition score and ride score of the sample road segment were tested for autocorrelation. As for the condition score data, I found the evidence of autocorrelation existence only when the segment contains the consecutive maximum values. The results of another analysis on the subsets without those consecutive sections did not show

autocorrelation. A reason for not detecting autocorrelation might be the property of subset data. Therefore, another road segment was tested with the same procedure, and it resulted in a subset without the maximum condition scores can have autocorrelation.

Based on the fact that condition score data have autocorrelation however it is mainly because of data structure-the consecutive maximum values, ride score was determined to be used as a substitute. Ride score data evidently showed autocorrelation. Two subsets were tested to verify if the length of data affect the autocorrelation test. The results indicated shorter subsets also have autocorrelation but evidence was not as strong as the case of using whole segment data.

Overall, it is not yet clear how prevalent autocorrelation is in pavement data. That is because of confounding results such as the effect of consecutive maximum in condition score data, and different results between ACF/PACF and Durbin-Watson test. However, intuitively, it is likely to have dependency in terms of pavement performance between neighboring sections. Also, the test using ride score data showed that the data definitely have AR(1) process.

Therefore, for the future research, when developing dynamic segmentation models AR(1) should be taken into account. On top of that, when AR(1) is considered in a model, ride score or other performance measure that have no consecutive maximum or minimum values presented should be used. For example, models using a Bayesian approach without AR(1) and with AR(1) can be developed and compared to determine which model fits better in performance measure data.

3.3 OFF-THE-SHELF CODES

There are several existing off-the-shelve packages for implementing various change point detection methods in R. Also, Hidden Markov Model (HMM) toolkits exist

in MATLAB. I tested some of them to see how the different methods work on road performance data in PMIS.

3.3.1 Change-point Detection Packages in R

Various packages have been developed to implement the change-point detection methods in R. The packages such as `changept`, `bcp`, `ecp`, `cpm`, `strucchange`, etc. were found throughout the search in the study. In this report, only three packages, including `changept`, `bcp`, and `ecp`, were used in data analysis.

3.3.1.1 R package: *changept*

The `changept` package (Killick, Eckley, and Haynes, 2016) provides multiple change-point search methods such as binary segmentation, segment neighborhood, and PELT. These methods are available both for changes in mean and/or variance by using assumption of either independent normal distribution or nonparametric cumulative sum (Killick and Eckley, 2014).

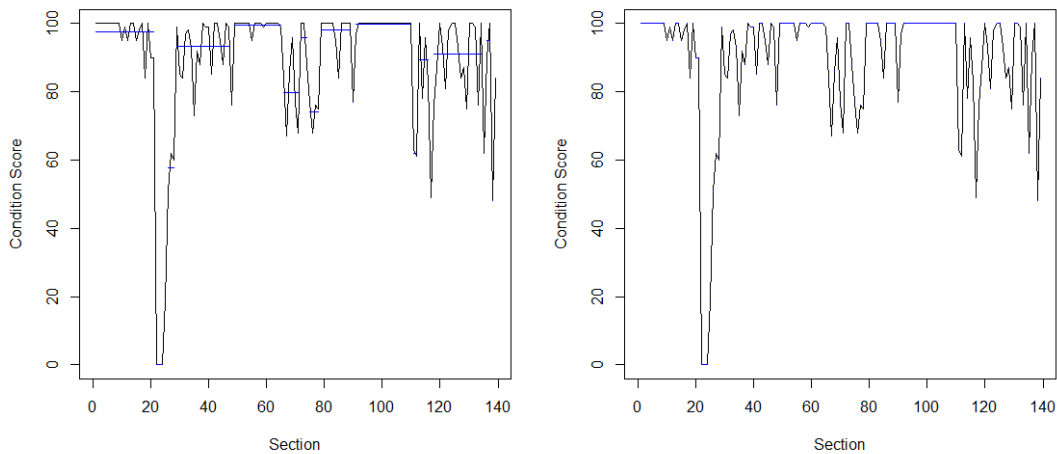


Figure 15: Resulting Plots of the `changept` Package in R.

In this package, multiple change points are able to be detected by minimizing sum of cost and penalty function. The penalty function prevents overfitting problem. For testing, function for detecting changes in mean was used with PELT algorithm. In Figure 15, a blue line indicates the mean value of a segment so the length of the blue line means the length of a segment. A default penalty value $\log(n)$ was used for the right plot, and adjusted penalty value $50 \cdot \log(n)$ was used for the left plot. On the right, almost all data points were detected as change points, which overfitting. By increasing the penalty, the number of change points decreased as shown on the left plot. However, the number of change points was sensitive to the penalty. Therefore, it seemed tricky to adjust the penalty value to obtain optimal number of change point by inspecting the plot.

3.3.1.2 R package: bcp

The *bcp* package is the implementation of Bayesian change point procedure by Barry and Hartigan (1993). A probability distribution is provided in the Bayesian procedure instead of specific locations of change points. This package provides both univariate and multivariate change point analysis by Markov chain Monte Carlo (C. Erdman and J. Emerson, 2007).

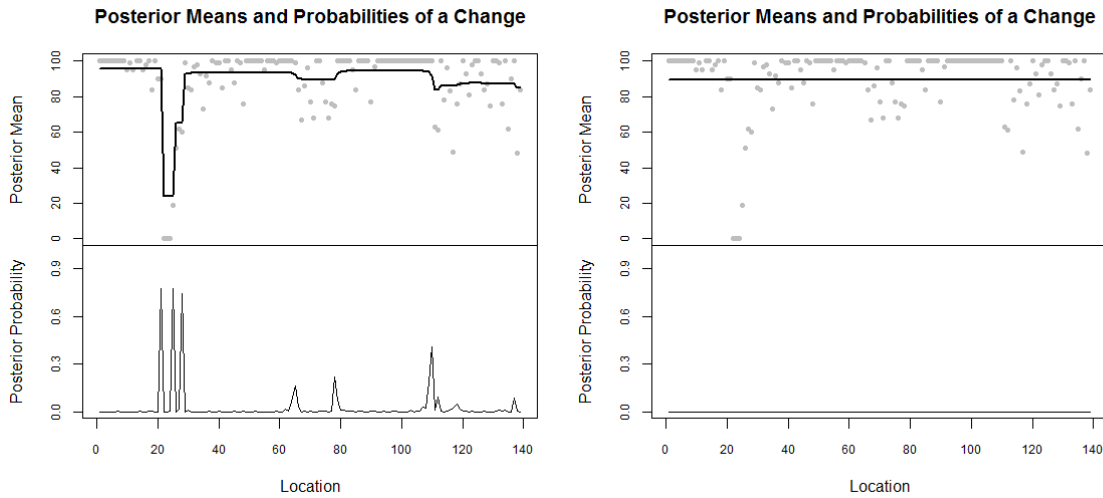


Figure 16: Resulting Plots of the bcp Package in R.

Figure 16 is the resulting plots of the bcp package when applying the test section data. The result shows the posterior mean of each partition after identifying change points. Also, the posterior probability is calculated so that change points can be chosen. In the bcp package, user can input optional values p_0 and w_0 . On the left, user input p_0 and w_0 were adjusted as $p_0=0.01$ and $w_0=0.01$ to detect the changes in means. On the right, default user inputs, both p_0 and w_0 equal to 0.2, were used, as a result no change point was identified. It seemed critical to choose appropriate user inputs for the parameters to obtain better segmentation result; however, the way to determine p_0 and w_0 is somewhat arbitrary because the result should be inspected visually to verify how well the algorithm partitions the section after choosing the user inputs.

3.3.1.3 R package: ecp

The ecp is an R package for nonparametric multiple change point analysis of multivariate data. Energy statistic was used to detect distributional change in a time series (James and Matteson, 2014)

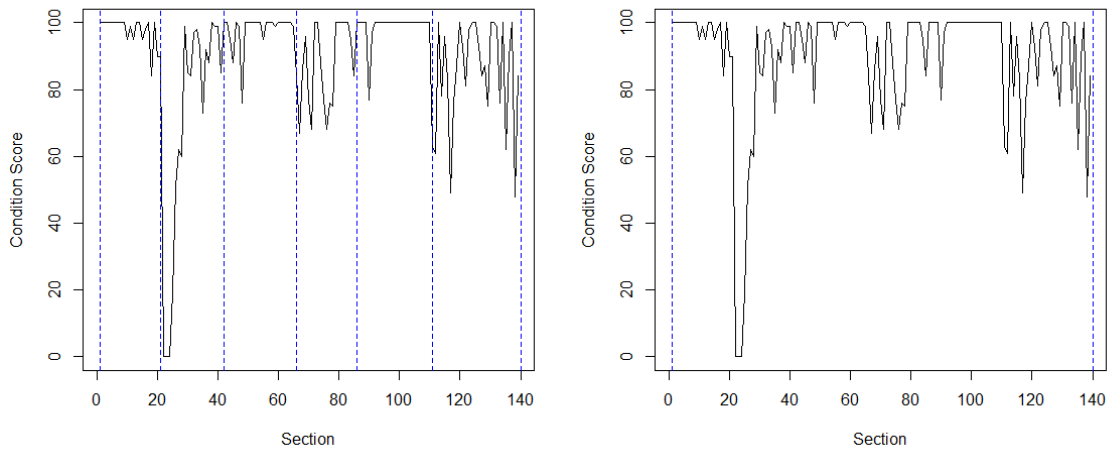


Figure 17: Resulting Plots of the ecp Package in R.

I chose to try the divisive algorithm which is recommended by the developer of the package. The algorithm sequentially identifies change points via a bisection algorithm. User inputs such as alpha, number of minimum observation, and maximum number of permutation were needed to determine. I used $\alpha = 1$, maximum number of permutation, $R = 500$, and varied number of minimum observation because other than this input are insensitive to the test data. Figure 17 shows the plots resulted from varying the number of minimum observation. The blue dashed line indicates the locations of change points for the sample data. On the left and right plots, the corresponding number of minimum observation was 20 and 30, respectively. According to a visual inspection, it seemed the ecp package reasonably detects the changes in distribution of each segment when user selects appropriate inputs (for instance min. observation = 20).

3.3.1.4 Discussion of Change-point Detection Packages in R

Throughout the trials of three off-the-shelf packages in R, it is found that various methods for change point detection were already implemented and R platform provides convenient solution for the change point detection problem.

However, each package has a few limitations. First, not all packages support multivariate data. That means segmentation is done by one variable if a univariate case. Ultimately, the goal of this study is developing a segmentation method to use multivariate data. For example, using both ride score and distress score simultaneously to detect change points. Therefore, these packages are not suitable to use as is.

On top of that, pavement performance data might have autocorrelation in nature. That is, each observation is not independent but has correlation so that the performance measure of current section has something to do with that of the next section. All packages that I tested in this study cannot take into account autocorrelation. The study of existence of autocorrelation will be discussed later in this chapter.

In addition to that, there was no feature for treating missing data. Because pavement management data usually have some portion of missing data due to undergoing construction or other reasons, it is important to take into account how to treat the missing data.

Lastly, a lack of evaluating the outcome of segmentation is problematic. The only way to check whether or not the result is good is visual inspection. The results of all packages were significantly affected by user inputs such as penalty and parameters. Therefore, it is critical to evaluate the result by finding appropriating user inputs. A developing method would have a feature to select these user inputs automatically. I suggest a way to evaluate the segmentation result in general using cross-validation later in this chapter.

3.3.2 Hidden Markov Model in MATLAB

3.3.2.1 Hidden Markov Model

A Hidden Markov Model (HMM) is defined as *a doubly stochastic process with an underlying stochastic process that is not observable, but can only be observed through another set of stochastic processes that produce the sequence of observed symbols* (Rabiner and Juang, 1986). Figure 18 illustrates the structure of HMM, Y_s are hidden states and X_s are observations. Two major assumptions are made by the model. First, the current state only depends on the previous state. Second, the observations are independent each other, it only depends on the current state.

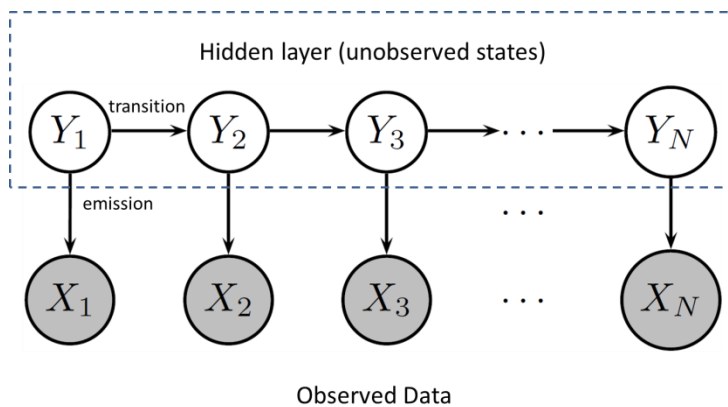


Figure 18: Graphical representation of HMM

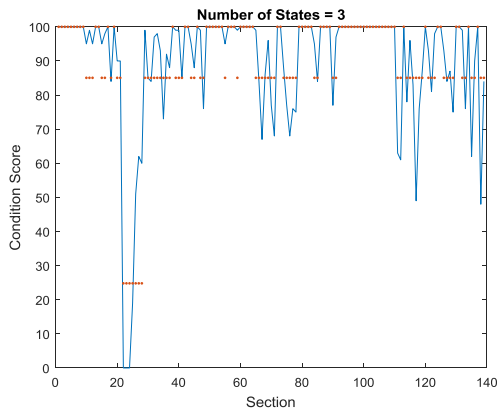
The HMM has three key problems such as classification, decoding and learning. The first problem is to compute the probability of the observed sequence produce by the model. The solution for this problem makes it possible to evaluate the model when there are several competing models. The forward and backward algorithm is the solution of the problem. The second problem is to find the most probable path given a set of observations. In other words, this solution attempts to reveal the hidden part of the model. One of the solutions for this problem is the Viterbi algorithm which results in the single

best state sequence based on an observation sequence. The third problem is how to optimize the model parameter for the given observed sequence. The expectation-maximization (EM) for HMM application, Baum Welch algorithm is the one of approaches to solve this problem (Blunsom, 2004). For more details, see references (rabinar, 1989; rabinar and Juang, 1986)

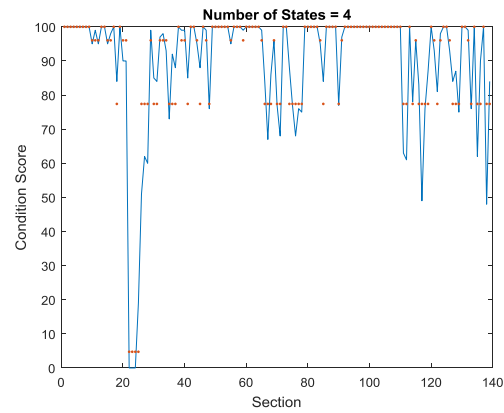
3.3.2.2 A toolkit for Hidden Markov Model in MATLAB

A MATLAB toolbox written by Kevin Murphy under MIT license (Murphy, 1998) was used to analyze the condition score data. Among provided features such as discrete HMM and mixed Gaussian, the mixed Gaussian HMM whose responses conditional on states are Gaussian was used because the condition score is not discrete but continuous.

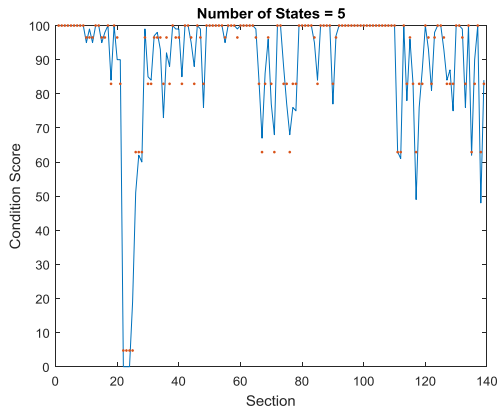
A procedure of the segmentation using HMM was conducted as follows. First, the number of states was defined. Then, parameters of HMM was estimated using the expectation-maximization (EM) algorithm-also known as Baum Welch algorithm (Baum et. al., 1970). The EM initialized parameters based on the number of states and Gaussian distribution and then computed the hidden state sequence based on the current parameters. Again, the parameters were re-estimated using the current hidden state sequence. This procedure iterated until convergence to a local maximum of the log-likelihood. As a next step, the most probable sequence is calculated by Viterbi algorithm (Viterbi, 1967) which computes the globally optimal state sequence using the estimated parameters.



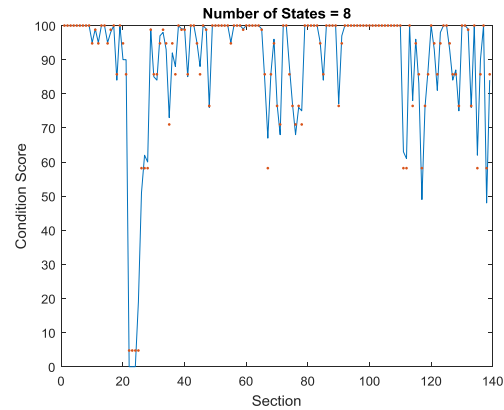
(a)



(b)



(c)



(d)

Figure 19: HMM results with varying the number of states equals to (a) 3; (b) 4; (c) 5; (d) 8

Figure 19 shows the results of HMM algorithm using various number of states. The number of states varied from 3 to 10; however, here I presented results of the number of states equals to 3, 4, 5 and 8. The blue line depicts condition score and the red dots are the mean values of each state. As the number of states increases, the number of segments increases. Therefore, when the number of states was 8, there were too many segments

delineated. On the other hand, the number of states 3 was not adequate for representing each state due to its high variance.

As this analysis was conducted to see the eligibility of HMM for segmentation, detailed adjustments, including a minimum section length, a minimum difference between states, etc., were not taken into account. Therefore, the results are not directly implemented to the segmentation application as is.

3.3.2.3 Discussion of Hidden Markov Model in MATLAB

As I mentioned previously, HMM was conducted to see if it can be used as a tool for the segmentation. The preliminary results showed that HMM is a promising mean for developing a sophisticated segmentation method due to its properties such as expandability and flexibility.

In order to develop a segmentation method using HMM, a few things should be taken into account. Firstly, it is important to note that the initialization is critical because the EM algorithm only finds a local optimum. For this reason, it was observed that several runs of the EM procedures with the same input gave different results due to the randomized initial parameters. Thus, a more refined way to determining the initial values need to be implemented. As an alternative, a method using Markov chain Monte Carlo (MCMC) would be used instead of optimization-based methods.

Also, some detailed adjustments, including a minimum section length, a minimum difference between states, etc., need to be set up. In the preliminary results, owing to the lack of those adjustments, the number of segments and the minimum section length were not able to be controlled. In addition, the number of states can be determined automatically using an optimization.

In this study, the most basic HMM with an assumption that observed variables are conditionally independent each other was used. For more realistic modeling to take into account spatial correlation of data, a HMM with the first order autoregressive, AR(1), whose observations are not independent anymore but have autocorrelations, would be a good alternative. Furthermore, not only univariate series data but also multivariate series can be analyzed using HMM so that multiple criteria simultaneously affect to the segmentation results. Additionally, a way to handle missing data is also a key issue to resolve.

3.4 IMPLEMENT A BAYES APPROACH

A Bayes approach (Thomas, 2001) was implemented to be used as a baseline method to compare the developing methods. Because there was no available off-the-shelf code for this approach, R was used to write a code. Among a series of methods suggested in Thomas' thesis, the method introduced in Paper I was implemented in this study to understand the basic concept of how the Bayesian approach used in the change point detection and expanded further to incorporate problems of change in variance, autocorrelation and multivariate series.

3.4.1 Data Used in Analysis

Ride score data for the same test section SH0046 was used instead of condition score data of that as the previous off-the-shelf package trials. The use of ride score was decided after testing autocorrelations in data, which will be discussed in the next chapter. One of the reasons why the ride score was chosen was that it was more appropriate to conduct the Bayesian approach with taking into account autocorrelations. Missing data were eliminated from the original data because the algorithm has no feature to handle the

missing values. As a result, total 135 sections were used in the ride score data. Figure 20 shows the plot of ride score and the histogram for the test section.

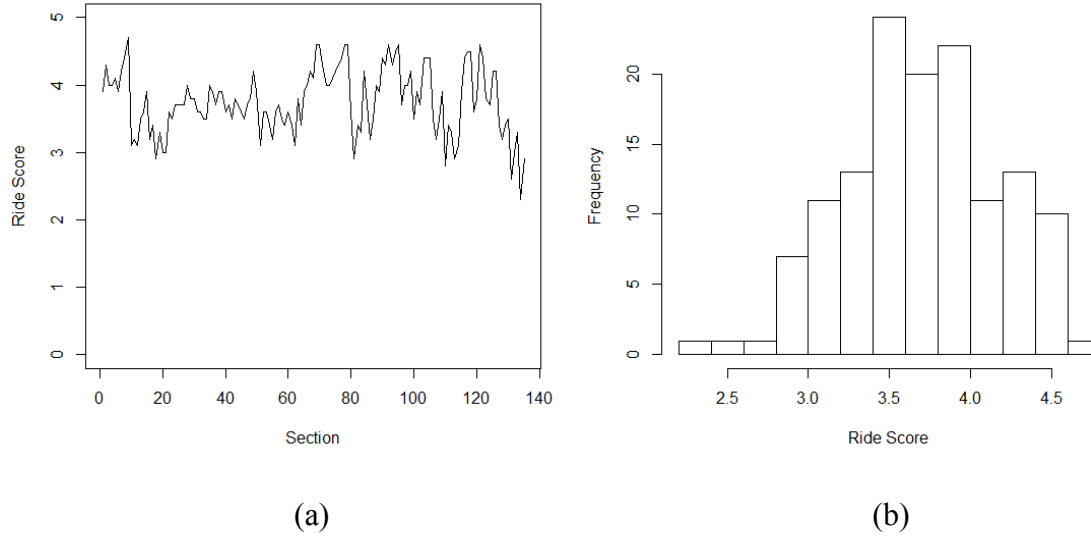


Figure 20: (a) Plot and (b) histogram of the ride score data

3.4.2 Data Analysis

3.4.2.1 Bayesian Approach with Unknown Common Variance

The basic concept of a Bayesian approach to detect a change point is already introduced in Chapter 2. Steps for detecting multiple change points using the Bayes approach were as follows. First, Bayes factors B_{r0} for $r = 1, \dots, 134$ were obtained by the equation for calculating the Bayes factor for M_r against M_0 ,

$$B_{r0} = \left(\frac{2}{3}\right)^{\frac{1}{2}} \left(\frac{n}{r(n-r)}\right)^{1/2} \left[\frac{S_{<n>}}{S_{<r>} + S_{<n-r>}}\right]^{n/2}$$

where, $S_{<n>} = \sum_{i=1}^n (x_i - \bar{x}_{<n>})^2$, $S_{<r>} = \sum_{i=1}^r (x_i - \bar{x}_{<r>})^2$, and $S_{<n-r>} = \sum_{i=1}^n (x_i - \bar{x}_{<n-r>})^2$

Then, calculate the Bayes factor for ‘change’ vs. ‘no change’ using,

$$B_{change} = \sum_{r=1}^{n-1} B_{r0} \frac{p(M_r)}{1 - p(M_0)}$$

The resulting $B_{change} = 362.9$, which is significant evidence of a change point, so the posterior probabilities for M_r could be calculated by,

$$p(M_r | x_1, \dots, x_n) = (r(n-r))^{-\frac{1}{2}} [S_{<r>} + S_{<n-r>}]^{-n/2} \times p(M_r)$$

Figure 21 shows the result of the first run which bisects the whole segment. At $r = 130$, the posterior probability was the maximum, that is, the first change point.

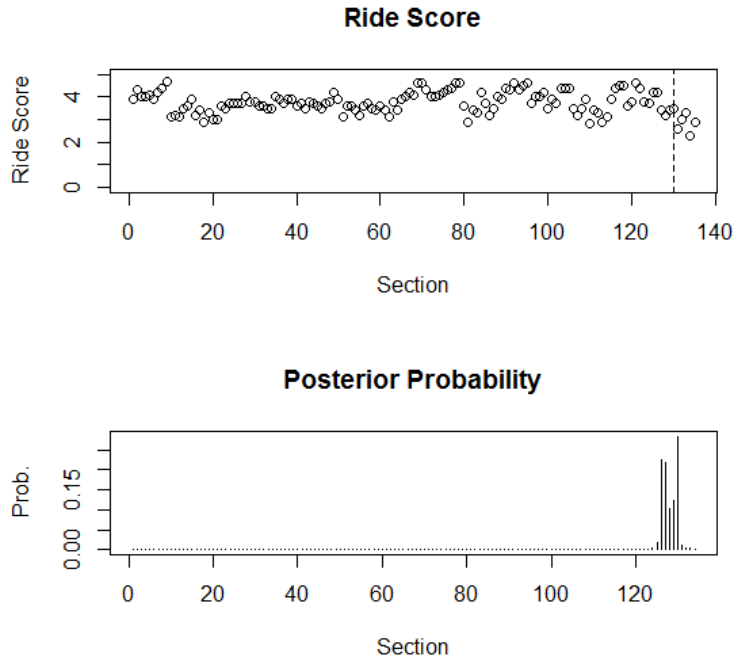


Figure 21: Result plot of the first run and posterior probabilities for a change point

In order to detect multiple change points, steps in the first run need to be repeated with dividend segments from the previous run. That is because the Bayesian algorithm is developed as at most one change-point (AMOC) which can detect only one change point at a time. For example, as the second run, two segments including sections at $r =$

1, ...,130 and $r = 131, \dots, 134$ should be analyzed. On top of that, to run the algorithm, a few user-specified inputs were needed to control the termination of iterations. When a Bayes factor for ‘change’ vs. ‘no change’ was less than 3, a decision was made that there is ‘no change’ for the corresponding segment. Also, a minimum number of section in a segment was determined as 2. Thus, the iteration will terminate when there is no segment left longer than 2 sections and having Bayes factor less than 3. Figure 22 displays the result of multiple change point detection by iterating the AMOC algorithm. In the result, there were eight segments were partitioned. For each segment, a mean and a range of two standard deviations from the mean were drawn as a solid red line and dashed lines, respectively.

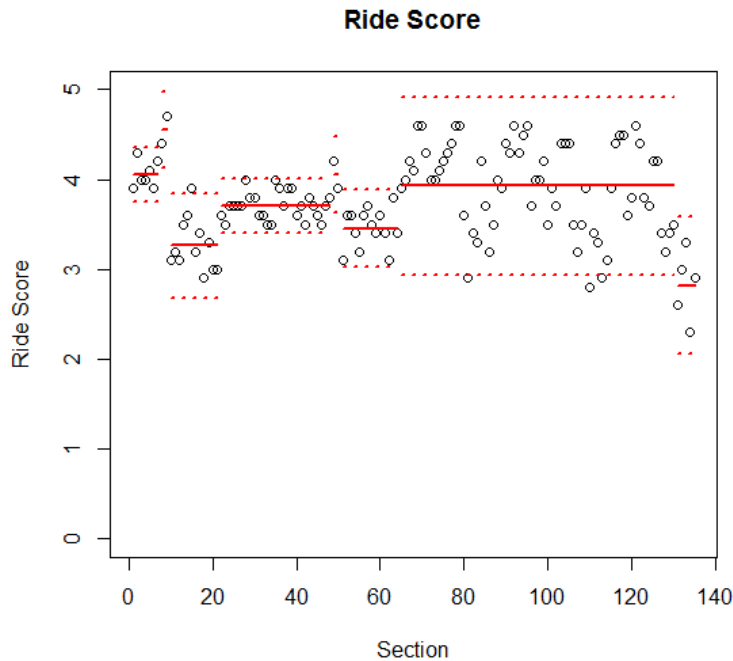


Figure 22: Result of detecting multiple change points using the Bayesian approach

The algorithm detects the change point one at a time due to it is designed as AMOC. Therefore, the iteration method is used to obtain multiple change points, and this ad-hoc

approach has limitations. One of the limitations could be seen in the result when we take a look at the 7th segment which has the longest length and the largest variance among all segments. Thus, it is a reasonable presumption that there could be more change points within the segment; however, the algorithm was not able to detect additional change points due to the limitation of AMOC.

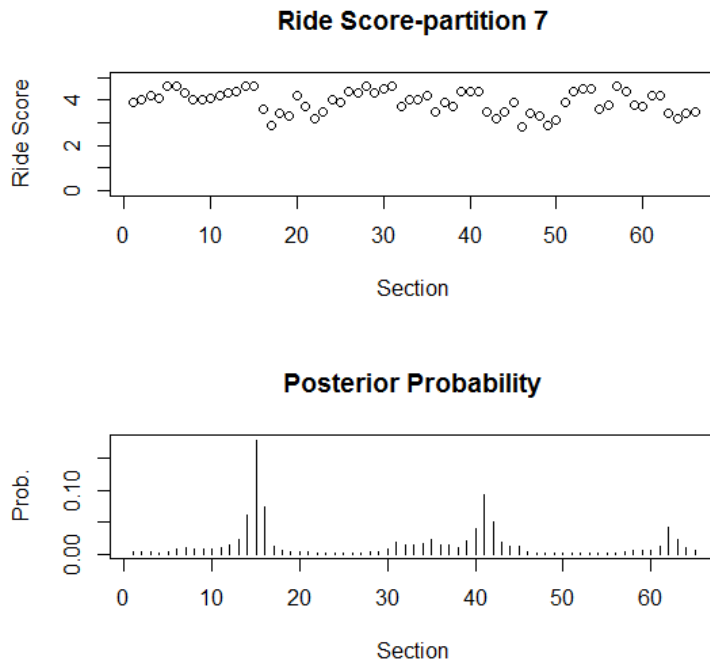


Figure 23: Result of the Bayesian algorithm illustrating a multi modal case using the 7th segment data

Figure 23 shows the result of the algorithm using the 7th segment data. Three peaks in the posterior probabilities can be seen in the graph. Although the posterior probabilities were not unimodal but multimodal, those modes were not able to be used to detect multiple change points because the posterior probabilities were calculated given a condition that there is at most one change point. The multimodal posterior case rather

hinders obtaining a significant Bayes factor. In this analysis, $B_{change} = 2.6$, this was less than the criteria 3 thus no more iteration progressed within this segment.

3.4.2.2 Bayesian Approach with First-order Autoregressive Process (AR(1))

With the same ride score data, a bit modified method was implemented. Instead of using the assumption that observations are independent, using AR(1) model to take into consideration the correlation between neighbored measures. Overall procedure was the same as non-AR(1) model; however, Bayes factors and the posterior probabilities were calculated by equations matrices involved in order to model the AR(1) form,

$$(z_t - \mu_t) = \varphi(z_{t-1} - \mu_{t-1}) + \varepsilon_t$$

where $|\varphi| < 1$, which means stationary process, and $\varepsilon_t \sim N(0, \sigma^2)$.

In order to run the algorithm φ needs to be estimated and input. $\hat{\varphi} = 0.61$ was estimated and used in the model, and the Bayes factor for ‘change’ vs. ‘no change’ resulted in $B_{change} = 2.9e-06$. That is, there was no significant evidence for a change in mean level of data. The result indicates that any existing variation is due to autocorrelation. This seems to disagree with a visual inspection, which suggests a segmentation around the locations detected by the non-AR(1) model. One reason for not detecting any mean level change may be because φ can vary along the segment. Therefore, a model which can detect the change in AR(1) together with detecting mean level seems more reasonable. In order to check the effect of φ on B_{change} , a plot of B_{change} vs. different φ values was drawn as shown in Figure 24. Only for a certain range of φ values has relatively significant Bayes factors, and that range is far from the estimated value used in the analysis.

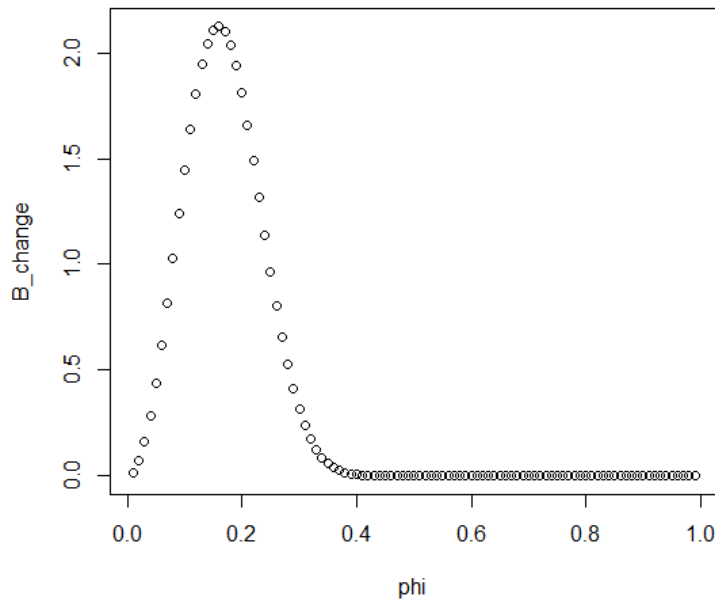


Figure 24: Bayes factor for ‘change’ vs. ‘no change’ vs. various values of φ

From the result above, although it is not meaningful to use other φ values, $\varphi = 0.2$ was tried in the model, which makes B_{change} value greater than 2, to specify what location has the maximum posterior probability. As a result, the maximum posterior probability was located at $r = 130$ which is consistent with the result from the approach without AR(1). Multiple change points were not able to be detected because of insignificant B_{change} value.

3.4.3 Discussion

A Bayesian approach to detect change points by Thomas (Thomas, 2001) were implemented using R. The ride score data were analyzed for both independent observation and autocorrelation algorithms.

It should be noted that this study examined only the first part of Thomas' thesis. Therefore the algorithm used in this analysis only detects changes in mean levels even though Thomas further developed methods to detect changes in variances in other papers. When it comes to AR(1) model, φ needed to be determined to run the model, however, Thomas' study presented a method to detect the change in φ .

Although the relatively crude model was used in the study, the problem of at most one change-point (AMOC) approach is still remarkable. The limitation of AMOC was verified through the analysis in the case of multimodal posterior probabilities. Thomas (2001) mentioned that the iterative method to find multiple change points is indefensible; however, a theoretically justified approach to detect multiple change points generates significant numbers of model, which is not feasible to compute (Thomas, 2001). As higher powered computers were appeared and more efficient algorithms were developed, it seems not impossible to develop a method to detect multiple change point at once. For instance, R package 'bcp' can detect multiple change points simultaneously using a Bayes approach even though the detailed algorithm was not investigated in the study.

The Bayesian approach is developed based on rigorous statistical background and also focused on road performance measures, this approach will be very useful as a baseline algorithm to evaluating various method being developed in near future.

Chapter 4: Conclusion and Future Study

This research is a preliminary study attempted to investigate the current techniques of segmentation for developing an improved method. Two main things have been studied for achieving the objective: (1) conducting the literature review, (2) testing the off-the-shelf tools and reproducing the Bayes approach. All things point to the conclusion that there are limitations for each approach and that there is good possibility that an improved method to overcome those limitations can be developed.

Unfortunately, developing an enhanced method could not be accomplished in this study; however, I would like to suggest directions for the future study. The most important point is that a developing approach should be based on a rigorous statistical method. In this sense, a HMM and a Bayesian approach are appropriate candidates not only because they are based on strong theoretical statistics but also because they are very flexible to resolve different issues. Once a method developed, for bench marking purpose, it seems good idea to compare several existing models such as CDA, R packages, random walk, and AR(1).

The main objective of future work is to develop a segmentation method capable of handling following issues have come across throughout this study.

Detecting multiple change points simultaneously

Majority of tested methods in this study delineate segment with identifying one change point at once and repeating algorithm to detect more changes using dividend segments in the previous run. Even though additional adjustments are suggested as a treatment, this approach does not result in optimal multiple change points. Either a HMM or a Bayesian approach has a capability to detect multiple change points simultaneously. Although that increases the difficulty of computational time, there would be solution

using the power of modern computer system and efficient optimizing algorithm. Therefore, the developing methods should have ability to identify the multiple homogeneous segments at once without losing optimality.

Taking into account AR(1)

The existence of autocorrelation in data is somewhat ambiguous even though it is examined by using testing methods such as Durbin-Watson, ACF and PACF. Thus, both models with and without AR(1) need to be developed and evaluated which model is more suitable to the application on a certain road performance measurement.

Model selection by cross validation

A method for comparing different approaches needs to be determined to evaluate segmentation methods. One of the ideas for comparing different models is to use a modified k-fold cross validation. As shown in Figure 25, observations in data are grouped in 5 fold. While holding the first group as a validation set, a model is estimated only using the rest 80% training set, then, repeat this procedure by shifting the validation set. Each model can be evaluated using the corresponding validation set by calculating log-likelihood.



Figure 25: Representation of 5-fold cross validation

Using R, the 5-fold cross validation method was implemented for the Bayesian approach. Figure 26 shows the result of segmentation using the first group as a validation

set and the rest group as a training set. Blue dots data in the validation set so that the log-likelihood could be calculated by the assumption of normal distribution in each segment.

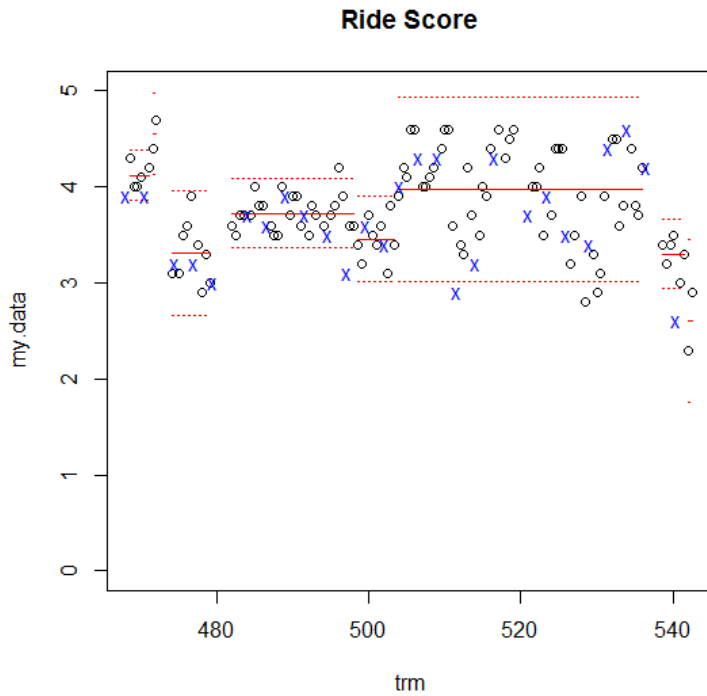


Figure 26: 5-fold cross validation for the Bayesian approach

Missing data treatment

When it comes to the road performance data collection, there are inevitable reasons for generating missing data. For example, road sections under construction for maintenance and rehabilitation cannot be inspected to collect the data. Thus, there should be a remedy for taking account into these missing data in the developing model. HMMs and most Bayesian models can resolve the problem. Hence, incorporating a treatment for missing data would be a significant improvement in the segmentation study.

References

- AASHTO. 1993. AASHTO Guide for Design of Pavement Structures. American Association of State Highway and Transportation Officials, Washington, DC.
- Acurio, J. R. M. 2014. Incorporating Risk and Uncertainty into Pavement Network Maintenance and Rehabilitation Budget Allocation Decisions. Ph.D. Dissertation. Texas A&M University.
- Barry, D., and Hartigan, J. A. 1993. A Bayesian Analysis for Change Point Problems. *Journal of the American Statistical Association*, 35(3), 309-319.
- Baum, L. E., Petrie, T., Soules, G., and Weiss, N. 1970. A maximization technique occurring in the statistical analysis of probabilistic functions of Markov chains. *Ann. Math. Statist.*, vol. 41, no. 1, pp. 164–171.
- Bennett, C. 2004. Sectioning of Road Data for Pavement Management. 6th International Conference on Managing Pavements, Brisbane, Australia, 1-11.
- Blunsom, P. 2004. Hidden Markov Models. Lecture Notes. Retrieved from <http://digital.cs.usu.edu/~cyan/CS7960/hmm-tutorial.pdf>
- Boroujerdian, M. A., Saffarzadeh, M., Yousefi, H., and Ghassemian, H. 2014. A Model to Identify High Crash Road Segments with the Dynamic Segmentation Method. *Accident Analysis & Prevention* 73 (December): 274–87.
- Breiman, L., Friedman, J. H., Olshen, R., and Stone, C. J. 1984. *Classification and Regression Tree* Wadsworth & Brooks/Cole Advanced Books & Software, Pacific California.
- Cafiso, S., and Graziano, A. D. 2012. Definition of Homogenous Sections in Road Pavement Measurements. *Procedia - Social and Behavioral Sciences*, SIIV-5th International Congress - Sustainability of Road Infrastructures 2012, 53 (October): 1069–79.
- Cuhadar, A., Shalaby, K., and Tasdoken, S. 2002. Automatic Segmentation of Pavement Condition Data Using Wavelet Transform. In *Canadian Conference on Electrical and Computer Engineering*, 2002. IEEE CCECE 2002, 2:1009–14 vol.2.
- Divinsky, M., Nesichi, S., and Livneh, M. 1997. Development of a Road Roughness Profile Delineation Procedure. *Journal of Testing and Evaluation*.
- Durbin, J. 1960. The fitting of time series models. *Rev. Inst. Int. Stat.*, v. 28, pp. 233–243.
- Durbin, J., and Watson, G. S. 1951. Testing for Serial Correlation in Least Squares Regression. II. *Biometrika* 38 (1/2): 159–77.

- Erdman, C., and Emerson, J. W. 2007. bcp: An R Package for Performing a Bayesian Analysis of Change Point Problems, *Journal of Statistical Software*, 23(3), 1-13.
- Haas, R. C. G., Hudson, W. R., and Zaniewski, J. P. 1994. *Modern Pavement Management*. Krieger Pub. Co, Malabar, FL.
- Jeffereys, H. 1998. *Theory of Probability*, 3rd edition. Clarendon Press, Oxford.
- Kennedy, J., Shalaby, A, and Cauwenberghe, R. V. 2000. Dynamic Segmentation of Pavement Surface Condition Data. 3rd Transportation Specialty Conference of the Canadian Society for Civil Engineering, London, Ontario.
- Killick, R., and Eckley, I. A. 2014. changepoint: An R Package for Changepoint Analysis. *Journal of Statistical Software*, 58(3), pp. 1–19.
- Killick, R., Haynes, K., and Eckley, I. A. 2016. changepoint: An R package for changepoint analysis. R package version 2.2.1, <http://CRAN.R-project.org/package=changepoint>.
- Kutner, M. H. 2005. *Applied Linear Statistical Models*. 5th ed. The McGraw-Hill/Irwin Series Operations and Decision Sciences. Boston: McGraw-Hill Irwin.
- Levinson, N. 1947. The Wiener RMS error criterion in filter design and prediction. *J. Math. Phys.*, v. 25, pp. 261–278.
- Misra, R., and Das. A. 2003. Identification of Homogeneous Sections from Road Data. *International Journal of Pavement Engineering*, 4:4, 229–233.
- Montgomery, D. C., Peck, E. A., and Vining, G. G. 2013. *Introduction to Linear Regression Analysis*. 5th ed. Chicester: Wiley.
- Murphy, K. 1998. Hidden Markov Model Toolbox for Matlab. Retrieved from <https://www.cs.ubc.ca/~murphyk/Software/HMM/hmm.html>
- Nicholas, J. A., and Matteson, D. S. 2014. ecp: An R Package for Nonparametric Multiple Change Point Analysis of Multivariate Data. *Journal of Statistical Software*, 62(7), 1-25.
- Ping, V. W., Yang, Z., Gan, L. and Dietrich, B. 1999. Development of Procedure for Automated Segmentation of Pavement Rut Data. *Transportation Research Record: Journal of the Transportation Research Board*, 1655, 65-73.
- Rabiner, L. R. 1989. A Tutorial on Hidden Markov Models and Selected Applications in Speech Recognition. *Proceedings of the IEEE*. 77, 257-286.
- Rabiner, L. R., and Juang, B. H. 1986. An Introduction to Hidden Markov Models. *IEEE Acousfics, Speech & Signal Processing, Magazine*, 3, 1-16.
- Raquel, P., and West, M. 2010. *Time Series: Modeling, Computation, and Inference*. Chapman & Hall/CRC Texts in Statistical Science Series. Boca Raton: CRC Press.

- Scullion, T., and Smith, R. 1997. TxDOT's Pavement Management Information System: Current Status and Future Directions. Report No. FHWA/TX-98/1420-S. Texas Transportation Institute, College Station, TX.
- Thomas, F. 2001. A Bayesian approach to retrospective detection of change-points in road surface measurements. PhD thesis, Dept. of Statistics, Stockholm Univ., Stockholm, Sweden.
- Thomas, F. 2001. Automated Road Segmentation Using a Bayesian Algorithm. *Journal of Transportation Engineering* 131 (8): 591–98.
- Thomas, F. 2003. Statistical Approach to Road Segmentation. *Journal of Transportation Engineering* 129 (3): 300–308.
- TxDOT. 2003. Managing Texas Pavements. Texas Department of Transportation Construction Division, Materials and Pavement Section.
- Viterbi, A. J. 1967. Error bounds for convolutional codes and an asymptotically optimum decoding algorithm. *IEEE Transactions on Information Theory* 13 (2): 260–269.
- Yang, C., Tsai, Y., and Wang, Z. 2009. Algorithm for Spatial Clustering of Pavement Segments. *Computer-Aided Civil and Infrastructure Engineering* 24 (2): 93–108.