

AN ABSTRACT OF THE THESIS OF

Michael M. Anderson for the degree of Master of Science in Computer Science
presented on June 13, 2013.

Title: Activity Detection on Free-Living Data Using Change-Point Detection

Abstract approved: _____

Weng-Keen Wong

Activity detection on time series sensor data is a rapidly emerging field with a lot of potential real-world application. In previous work researchers have unrealistically assumed that they know at which ticks in a time series the given subject stops doing one activity and starts doing another. In this thesis we explore the feasibility of segmenting realistic free-living time series data using techniques from change-point detection, and then classifying the predicted segments using standard supervised learning techniques. We compare this to the popular approach of splitting the time series into small windows of fixed-length, predicting on each window with a classifier, and then smoothing over the predictions with an HMM. We find that the HMM approach clearly outperforms the change-point detection approach, but that change-point detection may be promising given a modeling assumption that is appropriate to the data.

©Copyright by Michael M. Anderson
June 13, 2013
All Rights Reserved

Activity Detection on Free-Living Data Using Change-Point Detection

by

Michael M. Anderson

A THESIS

submitted to

Oregon State University

in partial fulfillment of
the requirements for the
degree of

Master of Science

Presented June 13, 2013
Commencement June 2013

Master of Science thesis of Michael M. Anderson presented on June 13, 2013.

APPROVED:

Major Professor, representing Computer Science

Director of the School of Electrical Engineering and Computer Science

Dean of the Graduate School

I understand that my thesis will become part of the permanent collection of Oregon State University libraries. My signature below authorizes release of my thesis to any reader upon request.

Michael M. Anderson, Author

ACKNOWLEDGEMENTS

I would like to acknowledge and thank my advisor, Dr. Weng-Keen Wong, for providing me with a clear and overarching vision of this project at all stages of its development, as well as his general expertise and helpfulness with all of the rough patches, sticking points, and unexpected problems that invariably accompany an endeavor of this magnitude. I would also like to thank the other members of my committee, Dr. Alan Fern, Dr. Raviv Raich, and Dr. Hector Vergara, for giving their time to help supervise the final stages of this project. Thanks go to all of my family and friends who encouraged and stood by me along the way. They are too numerous to name but they know who they are.

TABLE OF CONTENTS

	<u>Page</u>
1 Introduction	1
2 Related Work	3
3 Methodology	4
3.1 Overview	4
3.2 Datasets	4
3.2.1 Synthetic Data (OSU Hip)	4
3.2.2 LiME	6
3.3 Featurization	8
3.4 Base Classifiers	10
3.5 Top-Down Approach	11
3.5.1 Change-Point Detection	11
3.5.2 Experimental Setup	13
3.6 Bottom-Up Approach	15
3.6.1 HMMs	15
3.6.2 Experimental Setup	17
3.7 Performance Metrics	17
4 Results	20
4.1 Change-Point Detection	20
4.2 HMM	25
4.3 Discussion	29
5 Conclusion	31
Bibliography	31

LIST OF FIGURES

<u>Figure</u>	<u>Page</u>
3.1 Data Lifecycle of Change-Point Detection Experiments	5
3.2 Data Lifecycle of HMM Experiments	5
3.3 OSU Hip Activity Samples	7
3.4 LiME Day 1 Activity Samples	8
3.5 Reference and Test Data	12
3.6 Visual Interpretation of an HMM	16
3.7 An example time series with 20 ticks of data, from a dataset with three classes: A, B, and C. Four true class labels with corresponding windows are shown above the x-axis, and four predicted class labels with their corresponding windows are shown below the x-axis.	18
4.1 OSU Hip Results. Graphs are organized into rows by base classifier, and columns by evaluation metric. Change-point detection results were averaged over 30 splits into training, testing, and validation datasets, along with error bars showing a 95% confidence interval.	22
4.2 LiME Day 1 Results. Graphs are organized into rows by base classifier, and columns by evaluation metric. Change-point detection results were averaged over 30 splits into training, testing, and validation datasets, along with error bars showing a 95% confidence interval.	23
4.3 LiME Day 2 Results. Graphs are organized into rows by base classifier, and columns by evaluation metric. Change-point detection results were averaged over 30 splits into training, testing, and validation datasets, along with error bars showing a 95% confidence interval.	24
4.4 OSU Hip HMM Results. Graphs are organized into rows by base classifier, and columns by evaluation metric. HMM results were averaged over 10 splits into training (base classifier), validation, training (HMM), and testing datasets, along with error bars showing a 95% confidence interval.	26

LIST OF FIGURES (Continued)

Figure		Page
4.5	LiME Day 1 HMM Results. Graphs are organized into rows by base classifier, and columns by evaluation metric. HMM results were averaged over 10 splits into training (base classifier), validation, training (HMM), and testing datasets, along with error bars showing a 95% confidence interval.	27
4.6	LiME Day 2 HMM Results Graphs are organized into rows by base classifier, and columns by evaluation metric. HMM results were averaged over 10 splits into training (base classifier), validation, training (HMM), and testing datasets, along with error bars showing a 95% confidence interval.	28
4.7	Comparison of change-point detection and HMMs in terms of accuracy and detection time. For a given false positive rate, change-point detection algorithms split a time series into windows of a certain average size, which decreases as the false positive rate increases. From this it is possible to relate false positive rates from the CPD experiments to the window sizes of the HMM experiments. Here is the HMM data with false positive rates of {0.033, 0.028, 0.024, 0.021, 0.019, 0.017} that would corresponds to average window sizes of {10, 12, 14, 16, 18, 20}. Both graphs show decision tree results, and both show results from each of the three datasets clustered together over a single false positive rate.	30

LIST OF TABLES

<u>Table</u>		<u>Page</u>
3.1	Statistics used to convert time series windows into feature vectors (taken from [31]). T is the number of ticks in the given time series, x and v are individual axes of the time series data, and $x(i)$ is the i th data tick of x . .	9

Chapter 1: Introduction

One of the general goals of the field of artificial intelligence is to build computing devices that are “context-aware”, that act as more than just passive number-crunching machines that receive input data through very restrictive and wholly human-operated channels such as a keyboard or mouse. Context-aware devices are capable of using sensor data to understand the environment that they are situated in, such as the locations of nearby objects and how the objects are moving [1]. One subfield of context-aware computing that has been receiving considerable attention in recent years is activity detection. The goal of activity detection is to build computer plus sensor systems that are able to determine what activity a human subject is performing at any given moment.

Such systems have a variety of real-world applications. Researchers have been exploring the feasibility of using both wearable and non-wearable sensor systems to monitor the health of elderly patients that have or are at risk of developing degenerative physical and mental diseases [7]. The goal is to eventually build sensor-based monitoring systems that can aid doctors and family members in tracking the decline of subjects over time. Also, detection of an abnormal activity may indicate that a senior is undergoing a serious medical event such as a heart attack or slip-and-fall [30]. Another application of activity detection is to track the energy expenditure of subjects as they go through the course of their day. The traditional method of performing such tracking is with self-reporting by the subject of their activities. Wearable sensors offer an alternative approach that is not susceptible to misreporting due to bias, poor memory, or other confounding factors that a human reporter introduces into the system. One approach is to estimate the vigorousness or metabolic equivalent (MET) of an activity and calculate energy expenditure directly [23], while another is to attempt to predict the type of activity performed, and calculate energy expenditure using knowledge of how vigorous that activity is generally [28].

Activity detection generally assumes that sensor data will be represented as a time series, and that at any given moment in the time series the subject is performing one and only one type of activity. Thus the time series is thought of as being partitioned into

a number of non-overlapping intervals (*windows*), which are delimited by moments in time when the subject stopped performing one activity and started performing another. Previous work has treated activity detection as an offline problem, and has rarely considered performance metrics other than accuracy. In this work we are interested in the feasibility of partitioning and classifying a time series on free-living data in real time, so we will also evaluate our algorithms in terms of the amount of time required to detect that an activity change has occurred.

To predict changes in activity and partition the data, we used change-point detection, which is a field of statistics popular in control theory and other similar applications. We call this our top-down approach, because this method takes as input an initial time series and partitions it into smaller pieces using change-point detection. As an alternate approach we used the well known technique of partitioning the time series into small fixed-length non-overlapping windows, predicting the activity type of each window using a base classifier, treating that prediction as the observable state of an HMM, and finally solving the HMM for its hidden states. We call this our bottom-up approach, because we begin with small windows of fixed-length, and use an HMM to aggregate windows and smooth them together into larger activity intervals.

Chapter 2: Related Work

As mentioned in the previous chapter, sensor systems may consist of environmental or wearable devices. Some examples of environmental sensors that activity detection researchers have used to gather data are microphones [7], weight detection panels [18], cameras [6], and water usage detectors [7]. Researchers will generally place environmental sensors inside of a house, have subjects live in the house for a period of time, and attempt to predict for activity types like cooking, watching TV, etc.

Various wearable devices have been tried as well, such as RFID gloves [18], but the most popular wearable for activity detection purposes is the accelerometer. Besides being inexpensive, accelerometers tend to be small and lightweight, and so are fairly unobtrusive and user-friendly. Accelerometers also gather data at a high frequency, and as such may be used to collect a sizeable amount of data in a relatively short amount of time.

Whether or not an accelerometer will be discriminative for a set of activity types depends partially on where the accelerometer is worn on a subject's body. For example, an accelerometer worn on the ankle will be more discriminative for the activity of cycling than it would be if it was worn on the hip, and different types of arm movements will likely be discriminated only by an accelerometer worn on the arm. For this reason some researchers have opted to use multiple accelerometer systems to capture movement information from different parts of the body [4] [5]. However, this approach can be cumbersome for the wearer, so a single accelerometer is preferred when it is reasonable to assume that it will be discriminative over the relevant set of activities.

Activity sensor data tends to be noisy and not amenable to a deterministic or rule-based analysis, so activity types are typically modeled probabilistically, and activity detection is usually formulated as a supervised learning problem. The various common supervised learning algorithms are all familiar to the activity detection literature, though neural networks are especially popular, such as in [2] [21] [23]. More complicated meta-modeling approaches have also been tried, such as plurality voting with bagged, boosted, and stacked classifiers [15]; conditional random fields [29]; and HMMs [9] [14].

Chapter 3: Methodology

3.1 Overview

Each dataset that we tested consisted of multiple time series gathered from a number of different subjects, so to perform an experiment on a dataset we began by partitioning the set of time series into disjoint subsets of training, validation, and test data. Each individual time series was then partitioned into a set of disjoint windows, and each window was converted into its own feature vector. Once the dataset was featurized, the experiment could be treated as a typical classification problem. Base classifiers were built with the training set, and tuned (when necessary) on the validation set. To complete the change-point detection experiments (Figure 3.1), the tuned model was evaluated with the testing data. For the HMM experiments (Figure 3.2), the tuned model made predictions on a second training set, which were used to build an HMM metamodel. Finally, that metamodel was evaluated with the testing data. Following sections describe these processes in further detail.

3.2 Datasets

For our experiments we were interested in testing our algorithms on real-world free-living data. In the past, researchers have gathered activity data under unrealistic laboratory conditions, and by performing activities themselves instead of using independent subjects. Unfortunately there are not many such labeled free-living datasets available, so we only tested our algorithms on two datasets.

3.2.1 Synthetic Data (OSU Hip)

Our first dataset was collected by the Nutrition and Exercise Sciences department of Oregon State University, and has been used for previous activity detection research with the goal of designing a system to calculate and monitor subjects' energy expenditure [28] [32]. This dataset consisted of 91 time series collected over a 2-week period in a

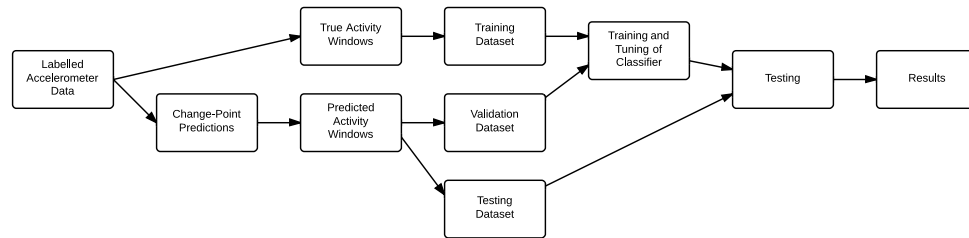


Figure 3.1: Data Lifecycle of Change-Point Detection Experiments

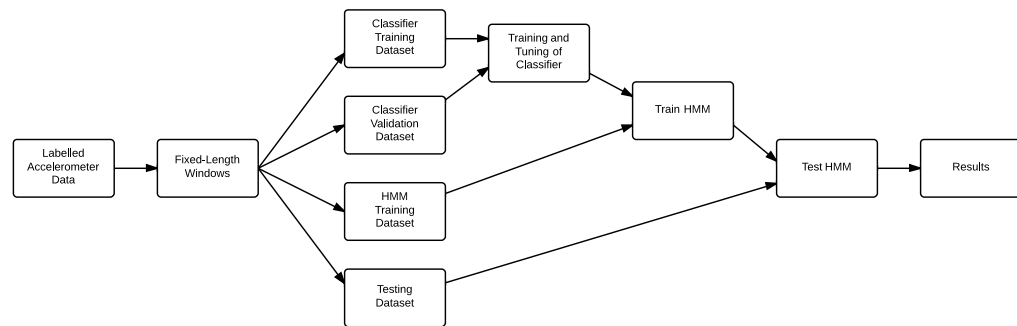


Figure 3.2: Data Lifecycle of HMM Experiments

laboratory environment, from 50 subjects who were children between the ages of 5 and 15. Subjects performed 12 different types of activities (as shown in Figure 3.3) over two separate visits, while an ActiGraph GT1M accelerometer worn on the hip collected triaxial acceleration data at a frequency of 30Hz.

Data was collected from two separate visits to the lab, where the subjects performed 6 activities per visit. Each subject performed the 12 activities in the same order. In the version of the dataset available to us, we had all 12 activities of 41 of the subjects, only the first 6 activities of an additional 5 subjects, and the last 6 activities of the remaining 4 subjects. Subjects were given breaks in between each activity and activities lasted 5-10 minutes, however, these unlabelled breaks were removed from the version that we used. Additionally, only two minutes of data were available for each subject, so time series consisted of six 120 second long activities that were concatenated together, qualifying the data as synthetic. Each of the 91 time series contained a total of $6 * 120 * 30 = 21600$ data ticks.

We determined that several of the activities were very similar and that it would be difficult to discriminate between them, so we combined some of them together to create a 7 class version of the data. Our classes were lying down, sitting (hand-writing, computer game), standing (laundry, sweeping, and catch), walking (comfortable, brisk and treadmill walking), dancing, running, and basketball.

3.2.2 LiME

This dataset consisted of 23 time series, each containing roughly 10 continuous days worth of data from an individual subject. It was collected by Helen Brown from the University of Cambridge, and Gemma Ryde from the University of Stirling, Scotland. Subjects wore an ActiGraph GT3X+ accelerometer during the entire period, which collected triaxial acceleration data at a frequency of 30Hz, as well as an activPal inclinometer on their thighs. The inclinometer provided what we considered the ground truth labels of the data by automatically delimitting and classifying intervals using the orientation of the subject's thigh at any given moment. It classified a horizontal orientation as lying down/sitting, a vertical orientation as standing, and a combination of the two as walking. Figure 3.4 shows samples of accelerometer data from the 3 activities.

This dataset was challenging to work with because of its size, as each individual time

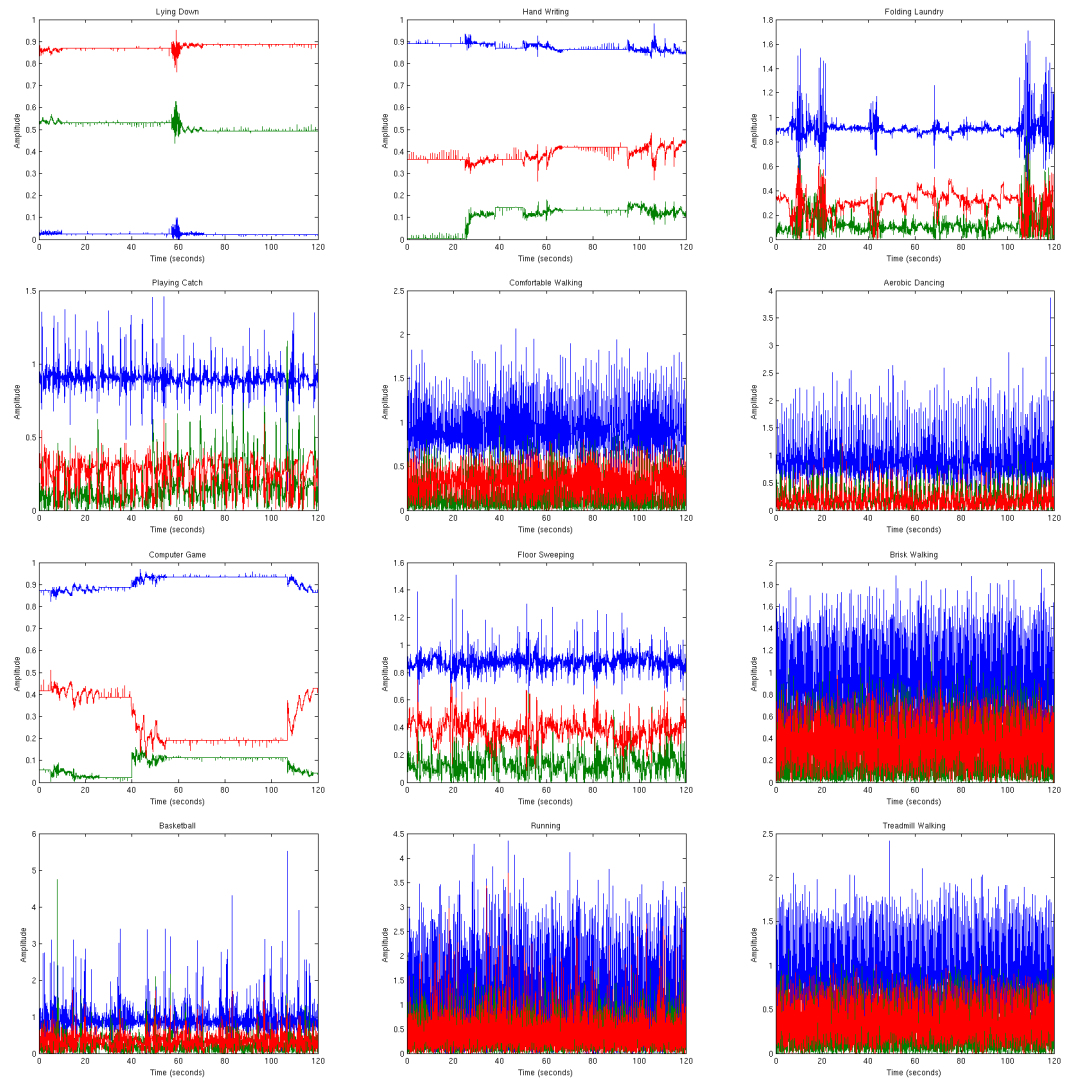


Figure 3.3: OSU Hip Activity Samples

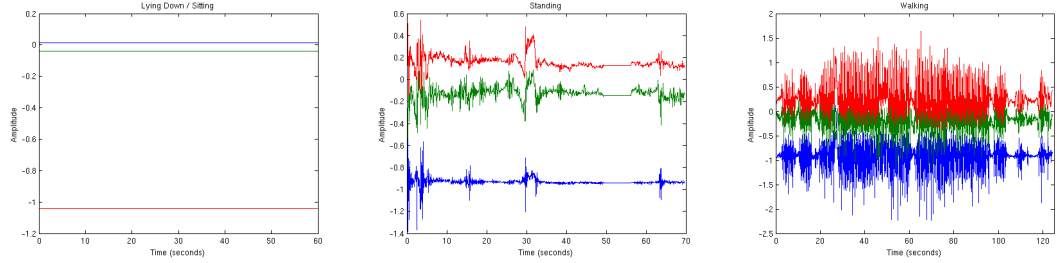


Figure 3.4: LiME Day 1 Activity Samples

series contained roughly 25 million ticks of data. To help alleviate this problem, we split each time series into individual days. We then treated the first full 24 hour period of data that began at midnight, from each subject, as one whole dataset (LiME Day 1), and the second such period as a separate dataset (LiME Day 2). We did not use any data from the remaining days.

In contrast to the OSU Hip dataset, LiME was not synthetic and activity lengths were variable. For LiME Day 1, the average activity length was 2977 with a standard deviation of 26428, and the median length was 348. The average number of activities per time series was 871. As would be expected, statistics for LiME Day 2 were comparable: the average activity length was 3108 with a standard deviation of 24099, and the median length was 348. The average number of activities per time series was 834. The medians were relatively small because many of the activities were short, while the mean and standard deviations were larger because a few of the activities were extremely long (e.g. when subjects were sleeping).

3.3 Featurization

To formulate our experiments as classification problems, we split each time series into a set of non-overlapping windows and represented each window as a feature vector. How we decided where one window ended (and where the next began) varied between experiments, and is described in sections 3.5 and 3.6. Our feature set was a large collection of statistics that have been shown to be discriminative for activity classification in previous research [10] [17] [23] [32]. In all we used 18 statistics that were uniaxial, i.e. were only a function of the data from a single axis of a given window, and one biaxial statistic.

The uniaxial statistics were applied to data from each axis separately, and the biaxial statistic was applied to data from each of the $C_2^3 = 3$ possible pairs of axes, for a total of $18 * 3 + 3 = 57$ features.

F1. Sum of values of a period of time: $\sum_{i=1}^T x(i)$.
F2. Mean: $\mu_x = \frac{1}{T} \sum_{i=1}^T x(i)$.
F3. Standard deviation: $\sigma_x = \sqrt{\frac{1}{T} \sum_{i=1}^T [x(i) - \mu_x]^2}$.
F4. Coefficients of variation: $\frac{\sigma_x}{\mu_x}$.
F5. Peak-to-peak amplitude: $\max\{x(1), \dots, x(T)\} - \min\{x(1), \dots, x(T)\}$.
F6-10. Percentiles: $10^{th}, 25^{th}, 50^{th}, 75^{th}, 90^{th}$.
F11. Interquartile range: difference between the 75^{th} and 25^{th} percentiles.
F12. Lag-one-autocorrelation: $\frac{\sum_{i=1}^{T-1} [x(i) - \mu_x][x(i+1) - \mu_x]}{\sum_{i=1}^{T-1} [x(i) - \mu_x]^2}$.
F13. Skewness: $\frac{\frac{1}{T} \sum_{i=1}^T [x(i) - \mu_x]^3}{(\frac{1}{T} \sum_{i=1}^T [x(i) - \mu_x]^2)^{\frac{3}{2}}}$, asymmetry of the signal probability distribution.
F14. Kurtosis: $\frac{\frac{1}{T} \sum_{i=1}^T [x(i) - \mu_x]^4}{(\frac{1}{T} \sum_{i=1}^T [x(i) - \mu_x]^2)^2} - 3$, peakedness of the signal probability distribution.
F15. Signal power: $\sum_{i=1}^T x(i)^2$.
F16. Log-energy: $\sum_{i=1}^T \log[x(i)^2]$.
F17. Peak intensity: number of signal peak appearances.
F18. Zero crossings: number of times the signal crosses its median.
F19. Correlation between each pair of axes: $\frac{\sum_{i=1}^T [x(i) - \mu_x][v(i) - \mu_v]}{\sqrt{\sum_{i=1}^T [x(i) - \mu_x]^2 \sum_{j=1}^T [v(j) - \mu_v]^2}}$.

Table 3.1: Statistics used to convert time series windows into feature vectors (taken from [31]). T is the number of ticks in the given time series, x and v are individual axes of the time series data, and $x(i)$ is the i th data tick of x .

One discriminative characteristic of an activity is its overall vigorousness. The sum [F1] and the sample mean [F2] both act as simple and obvious ways of measuring vigorousness, as more intense activities will tend to involve higher rates of acceleration during movement. We also used the 10th [F6], 25th [F7], 50th [F8], 75th [F9], and 90th [F10] percentiles of the data, as well as signal power [F15] and log energy [F16] as supplemental measures of overall activity intensity.

Another characteristic of an activity is how much it varies in intensity. The sample standard deviation [F3], coefficient of variation [F4], peak-to-peak amplitude [F5], number of zero crossings [F18], as well as the interquartile range [F11] were useful for discriminating between activities with a consistent level of intensity (low variance, etc.) and activities that were more rhythmic or staccato in intensity (high variance, etc.).

Lag-one-autocorrelation [F12], Skewness [F13], kurtosis [F14], and peak intensity [F17] were useful for discriminating between activities that tend to be similar in their overall intensity and variation in intensity, but that showed other types of difference in shape. Skewness indicates whether the data is more concentrated above or below its mean. Kurtosis indicates that the data is concentrated near its mean or conversely that it is fat-tailed. Lag-one-autocorrelation is a measure of the general relationship between data ticks and their immediate neighbors in time. Peak intensity is the number of times that the data repeatedly reached its maximum value.

Finally we looked at a single bimodal statistic across each pair of axes, the correlation coefficient [F19], which discriminates between activities where acceleration values in one axis are predictive of acceleration values in another axis, versus activities where that is not the case.

3.4 Base Classifiers

We tested 3 types of classification models on the featurized versions of our data: decision trees, support vector machines, and neural networks. We used R for our experiments, and used the R libraries ‘rpart’ [27], ‘e1071’ [12], and ‘nnet’ [16] to build our decision tree, svm, and neural net models, respectively. Default rpart values were used for the decision tree experiments, hence there were no parameters tuned and the validation set was ignored. For the neural net experiments, the maximum number of iterations was set to 100000, and the maximum number of weights was set to 1000000.

For the OSU Hip experiments we tuned the cost parameter c of the svm on the validation set with 6 values: $\{0.01, 0.1, 1, 10, 100, 1000\}$. The single-layer feed-forward neural network took two tuning parameters, and we tuned with each element of the set $N \times W$, where $N = \{1, 2, \dots, 30\}$ was the numbers of nodes in the hidden layer, and $W = \{0, 0.5, 1\}$ was the weight decay parameters.

Since the LiME datasets were an order of magnitude larger, we tuned them slightly differently because of time constraints. Setting the c parameter to 1000 proved to be very computationally expensive for the svm model, so we tuned c from the values $\{0.01, 0.1, 1, 10, 100\}$. Running $30 * 3 = 90$ tuning experiments for the neural networks was also prohibitively expensive, so we drew from $N \times W = \{5, 10, 15\} \times \{0, 0.5, 1\}$.

3.5 Top-Down Approach

3.5.1 Change-Point Detection

For this approach, the data was split into non-overlapping segments for featurization using techniques from the statistical field of change-point detection. Change-point detection has found application in many problem domains that require analysis of time series data from dynamic systems, including failure detection [3], quick detection of attacks on computer networks [25], and monitoring of heartbeat fluctuations during sleep [22]. Change-point detection techniques assume that each tick of a time series is a draw from some process, but that the process may suddenly change as time passes. The goal is to predict when these changes have occurred. A *score* is generated for each time tick, and if the score is above a given threshold a change is predicted to have occurred between that tick and its immediate predecessor. To generate a score at a time tick, a window of data that immediately preceeds it (the *reference data*) is compared to it along with a window of data that immediately follows it (the *test data*), as shown in Figure 3.5.

Model-based approaches to change-point detection assume that each tick in a time series is a draw from some underlying probability distribution. Scores are generated by estimating the distribution of the reference data and the test data, and then by calculating the likelihood that the two distributions are different. Where it is reasonable to assume that the data belongs to a particular family of distributions then parametric estimation methods have been employed [26]. If no such modeling assumptions are reasonable



Figure 3.5: Reference and Test Data

then non-parametric methods have also been found to be viable [11]. Distance-based approaches such as Singular Spectrum Analysis generate scores through other metrics of dissimilarity or difference between the reference data and the test data [13]. Notationally, we say that for each tick i in a time series:

$$s(i) = D(X_r(i), X_t(i))$$

Where $s(i)$ is the score of the i th tick, $X_r(i)$ is the reference data associated with the i th tick, $X_t(i)$ is the test data associated with the i th tick, and $D(A, B)$ is a function that computes the dissimilarity between a matrix of data A and matrix of data B , and varies between change-point algorithms. Note that for a given algorithm it may not be possible to generate scores right at the very beginning of the time series (insufficient reference data) or at the very end of a time series (insufficient test data).

3.5.2 Experimental Setup

There are many different modeling assumptions and associated algorithms for generating change-point detection scores, and one simple baseline approach that we wanted to test was the Shewhart Control Chart [20]. This approach assumes that the reference data is drawn from a multivariate normal distribution, and that scores are calculated by the Mahalanobis distance of the target time tick from the estimated multivariate normal:

$$s(i) = \sqrt{[\bar{X}_r(i) - X(i)]^T \frac{1}{S_r(i)} [\bar{X}_r(i) - X(i)]}$$

where $\bar{X}_r(i)$ is the sample mean of the reference data, $S_r(i)$ is the sample covariance matrix of the reference data, and $X(i) = X_t(i)$ is the i th data point.

We were also interested in testing the performance of a newer and more sophisticated change-point detection algorithm: the Kullback-Leibler Importance Estimation Procedure (KLIEP), introduced by Kawahara and Sugiyama [8] [24]. This approach generates scores using the Kullback-Leibler (KL) divergence between the reference data and the test data. One method of doing this is to estimate the density of the reference distribution and test distribution separately, and then compare them using a likelihood ratio (known in the change-point detection literature as *importance*). Instead, KLIEP estimates the importance directly using a non-parametric model.

Let the estimate of the importance \hat{R} be represented by this model:

$$\hat{R} = \frac{p_t}{\hat{p}_r} = \sum_{i=1}^{T_t} \alpha_j K_G(X_{rt}, X_t(j))$$

Where p_r and is the probability density of the reference data, p_t is the probability density of the test data, T_t is the number of ticks in the test window, α is a vector of model parameters to solve for, X_{rt} is the concatenation of the reference and the test data, $X_t(i)$ is the i th element of the test data, and $K_G[A, B]$ is the Gaussian kernel with width σ :

$$K_G[A, B] = \exp\left(-\frac{\|A - B\|^2}{2\sigma^2}\right)$$

Now solve for α so that the empirical KL divergence between \hat{p}_t and $p_t = \hat{R}p_r$ is

minimized, which is equivalent to the following convex optimization problem:

$$\begin{cases} \max_{\alpha} & \sum_{j=1}^{T_t} \log \left(\sum_{k=1}^{T_t} \alpha_k K_G[X_t(j), X_t(k)] \right) \\ \text{s.t.} & \frac{1}{T_r} \sum_{j=1}^{T_r} \sum_{k=1}^{T_t} \alpha_k K_G[X_r(j), X_t(k)] = 1 \\ & \text{and } \alpha_1 \dots \alpha_{T_t} \geq 1 \end{cases}$$

Finally, the scores that we wish to generate are just the estimate of the importance given by the solution to the complex optimization problem, i.e. $s(i) = \hat{R}(i)$.

Since this approach uses a Gaussian kernel, it requires the selection of a kernel width σ for each time tick. We used an implementation of KLIEP that is available at Sugiyama's website, which included a cross-validation procedure for the value of σ . The CV procedure chooses a number of disjoint splits of the test data along with a number of different candidate σ 's, and runs KLIEP with each combination of split and candidate σ . Then it chooses the candidate σ that, on the average across all of the splits, maximizes the KL divergence (the \max_{α} equation above) the most.

For the OSU Hip dataset, we used this CV procedure to choose the kernel width at each individual time tick. This computationally intensive approach was impractical for the UQ dataset because it is orders of magnitude larger, so instead of running it on every tick of that data, we ran the CV procedure on a number of random ticks drawn from the data. From this we were able to empirically identify 0.01 as a plausible σ , and so fixed σ at that value for our experiments on that dataset.

Our selection of reference and test window sizes were informed by two considerations: first that the window sizes contain enough data to accurately model the reference and test distributions, and second that the window sizes were small enough to detect activity changes at realtime speeds. Previous research [32] found that a reference window size of 10 seconds contained just enough information to discriminate well between OSU Hip activities. Since this window size worked well in previous experiments, and since the activities of the UQ dataset were comparable in average length, we decided to fix our reference window size at 10 seconds for both datasets. Because we were interested in minimizing detection time, and because 1 second was the smallest window that we felt could provide some information about an activity, we fixed our test window size at 1

second for both datasets.

Once change-point detection scores were generated, we tested a number of threshold values that determined which scores were high enough to be considered a predicted change-point. Threshold values were chosen by considering the false positive rates of change prediction for the change-point detection algorithms. A smaller false positive rate corresponded to a higher and more conservative threshold, which split the time series into fewer segments for featurization. A larger false positive rate corresponded to a lower threshold, which split the time series into more segments.

3.6 Bottom-Up Approach

3.6.1 HMMs

Once we had created our methodology for splitting up time series using change-point detection, we decided to test it against a more standard, baseline technique. For this approach, we used the Hidden Markov Model to take advantage of the sequential nature of our data. An HMM is a temporal graphical model that contains a set of hidden states $H = \{H_0, H_1, \dots, H_w\}$ as well as a set of observed states $O = \{O_1, O_2, \dots, O_w\}$ (Figure 3.6). An index of either type of state represents a point in time, such that if there exists two indices i and j where $i < j$, i is thought of as having happened before j . Each hidden state in the model has one of a discrete set of values associated with it drawn from $X = \{X_1, X_2, \dots, X_\ell\}$, and each observed state has one of a discrete set of values associated with it drawn from $Y = \{Y_1, Y_2, \dots, Y_m\}$. The values of the hidden states are unknown, and the values of the observable states are known. It is also assumed that, as indicated by Figure 3.6, the value of an observable state O_i is dependent on the corresponding hidden state H_i , and that the value of any hidden state H_i is dependent only on its immediate predecessor H_{i-1} .

Furthermore, the dependencies between hidden states and their followers are assumed to be described by a stationary stochastic process known as the *transition model*, $T : X^2 \rightarrow [0, 1]$. The dependencies between a hidden state and its adjacent observable state are assumed to be part of a separate but also stationary stochastic process known as the *observation model*, $S : X \times Y \rightarrow [0, 1]$. In other words, both models can be thought of as a function of two values, that output the probability of a change from

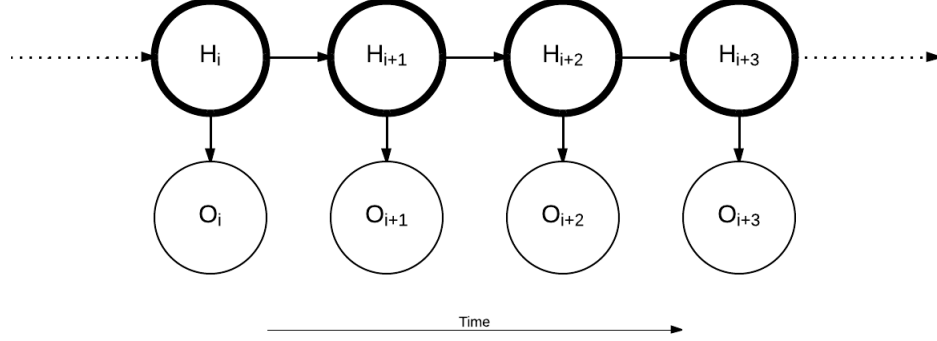


Figure 3.6: Visual Interpretation of an HMM

the first value to the second via a dependency arc in the HMM. The usual approach to estimating T and S given only O is a flavor of expectation maximization known as the Baum-Welch algorithm, which is useful for finding \hat{T} and \hat{S} that are locally maximally likely. However, suppose a *training HMM* $\langle H_{tr}, O_{tr} \rangle$ with the same model parameters is given, and the values of all of its hidden states as well as its observable states are known. T and S are then approximated by the following global maximum likelihood estimators:

$$\hat{T}(X_i, X_j) = \frac{|\{H_k \in H_{tr} \mid H_k = X_i \text{ and } H_{k+1} = X_j\}|}{|\{H_k \in H_{tr} \mid H_k = X_i\}|}$$

$$\hat{S}(X_i, Y_j) = \frac{|\{H_k \in H \mid H_k = X_i \text{ and } O_k = Y_j\}|}{|\{H_k \in H \mid H_k = X_i\}|}$$

Finally, if all of the values of O are known, and we are given a \hat{T} and an \hat{S} estimated from a training HMM, then the goal we are interested in is to use that information (along with the model assumptions) to find the most likely values for each state in H . There exists a polynomial-time dynamic programming solution to this problem known as the Viterbi algorithm. [19]

3.6.2 Experimental Setup

For our experiments we began by splitting each time series into small non-overlapping windows. Within a given experiment the window size was fixed, but across different experiments we tested window sizes of length $\{1, 2, \dots, 20\}$. Once the time series were split they were featurized. Classification models were built with training data, and tuned (in the case of the svm and neural net models) using validation data, in the same way that has already been described previously in this chapter.

Unlike the change-point detection experiments, this experiment required that the data be split into 4 equal parts: training (classifier), validation, training (HMM), and testing) rather than 3. Here we formulated the problem of making predictions on the testing set in terms of an HMM, first by treating the second training set as a training HMM. Each window of the second training set was treated as a time index $1, 2, \dots, n$. In our datasets we let H be the ground truth activity classes of the windows, and O be the predicted activity classes of the windows. We used the procedure above to calculate \hat{T} and \hat{S} , and assumed that these estimates held for the testing set as well as the second training set. We then used the tuned base classifier to predict on the testing set, giving us O . Finally, we used O , \hat{T} , and \hat{S} to run the Viterbi algorithm on the testing set and predict the ground truth activity classes H .

3.7 Performance Metrics

To measure the performance of our classification algorithms we used two metrics. Accuracy is defined as the number of ticks that an algorithm correctly classifies in a time series, over the total number of ticks in the time series. Since we were also interested in our algorithms' feasibility for activity classification in real time, we used detection time as a second metric. Detection time is defined as the average amount of time required for an algorithm to begin correctly classifying data after a true activity change has occurred.

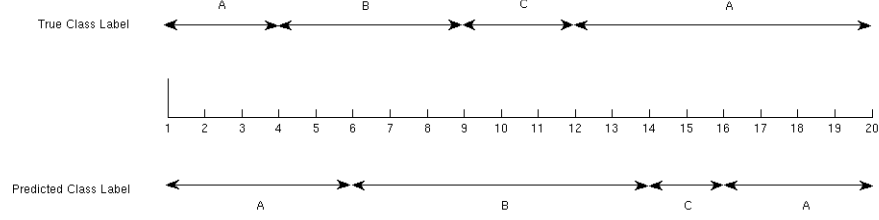


Figure 3.7: An example time series with 20 ticks of data, from a dataset with three classes: A, B, and C. Four true class labels with corresponding windows are shown above the x-axis, and four predicted class labels with their corresponding windows are shown below the x-axis.

To illustrate the metrics, an example section of a time series is shown in Figure 3.7. Accuracy can be computed by counting the number of correctly predicted ticks for each true window separately, summing the counts, and dividing by the total number of ticks.

$$\text{Accuracy} = \frac{\text{Correct}_{A1} + \text{Correct}_B + \text{Correct}_C + \text{Correct}_{A2}}{(\text{Total number of ticks})} =$$

$$\frac{3 + 3 + 0 + 4}{20} = 50\%$$

Detection time is computed by counting how many ticks are required after a true window has begun for the prediction algorithm to start correctly predicting the class. The counts are summed, and then divided by the number of true activities. Over the true window beginning at tick 1 (class A), the algorithm predicts A immediately, so the detection time for that window is 0. Over the second true window beginning at tick 4 (class B), the algorithm does not start predicting B until tick 6 so the detection time for that window is $6 - 4 = 2$. Over the third true window starting at tick 9 (class C), the algorithm never predicts C, so the detection time for that window is 3, the full length of the window. Over the fourth true window starting at tick 12 (class A), the algorithm starts to predict C at tick 14, but does not predict A until tick 16, so the detection time for that window is $16 - 12 = 4$. Thus the detection time for the entire time series is:

$$\frac{(0 + 2 + 3 + 4) \text{ ticks}}{4 \text{ windows}} = \frac{9}{4} \text{ ticks per window}$$

In the change-point detection experiments accuracy and detection time were averaged over 30 random splits of the given dataset into training, validation, and testing sets. Because the HMM experiments were more computationally expensive, accuracy and detection time were averaged over 10 random splits of the given dataset into training (base classifier), validation, training (HMM), and testing sets.

Chapter 4: Results

4.1 Change-Point Detection

Results for our change-point detection experiments are given in Figures 4.1-4.3. We hypothesized that the performance of the change-point detection algorithms would depend heavily on the threshold level for change prediction. This was tested by varying the average number of times per second that the algorithms falsely predicted a change. A large number of such false positive rates per second were tested, but for the sake of brevity only a representative sample of $\{0.005, 0.01, 0.05, 0.1\}$ are shown here.

In the OSU Hip experiments, control charts outperformed KLIEP in terms of detection time (Figures 4.1.2, 4.1.4, 4.1.6), while the accuracy results (Figures 4.1.1, 4.1.3, 4.1.5) were mixed. Except when predicted windows are large enough to span across multiple true activities, it is generally expected that accuracy will decrease as false positive rate increases because small windows contain less information and are less discriminative than larger windows. This behavior is seen in the control chart accuracy results (grey bars in Figures 4.1.1, 4.1.3, 4.1.5), but not in the KLIEP accuracy results (white bars in Figures 4.1.1, 4.1.3, 4.1.5). Follow-up experiments showed that KLIEP peaks in accuracy for false positive rates between 0.2 and 0.3 for all three classifiers. KLIEP seemed to perform best on this dataset when it was given many opportunities to predict changes.

Further investigation indicated that across the OSU Hip dataset the KLIEP algorithm was unable to detect many of the different activity changes without a very low score threshold value (and a very high false positive rates). Some qualitative plotting of the OSU Hip data showed that most of its activities have accelerometer amplitude values that strongly resemble draws from a multivariate normal distribution. Since control charts assume that the data is drawn from a distribution that is a member of that family, it is logical that control charts would outperform algorithms with different modeling assumptions on OSU Hip.

In the LiME experiments, KLIEP outperformed control charts in terms of accuracy across the board, and control charts outperformed KLIEP in terms of detection time

across the board. This suggests that in general control charts correctly detected true changes more quickly, but that after a correct change prediction it was more likely to make an incorrect change prediction.

In a few cases (Figures 4.1.2, 4.1.6, 4.2.6) the detection time did not decrease as the false positive rate increased. On the face of it this would seem to be a non-sequitur, but this only happened in cases when accuracy also decreased (Figures 4.1.1, 4.1.5, 4.2.5). Smaller window sizes tend to be correlated with decreased detection times, but it is possible that predicting with smaller windows, if they happen to contain an insufficient amount of discriminative data, can actually increase the time required for the classifier to start correctly predicting the ground-truth activity. Additionally, the given increases in detection time were small and within confidence bounds.

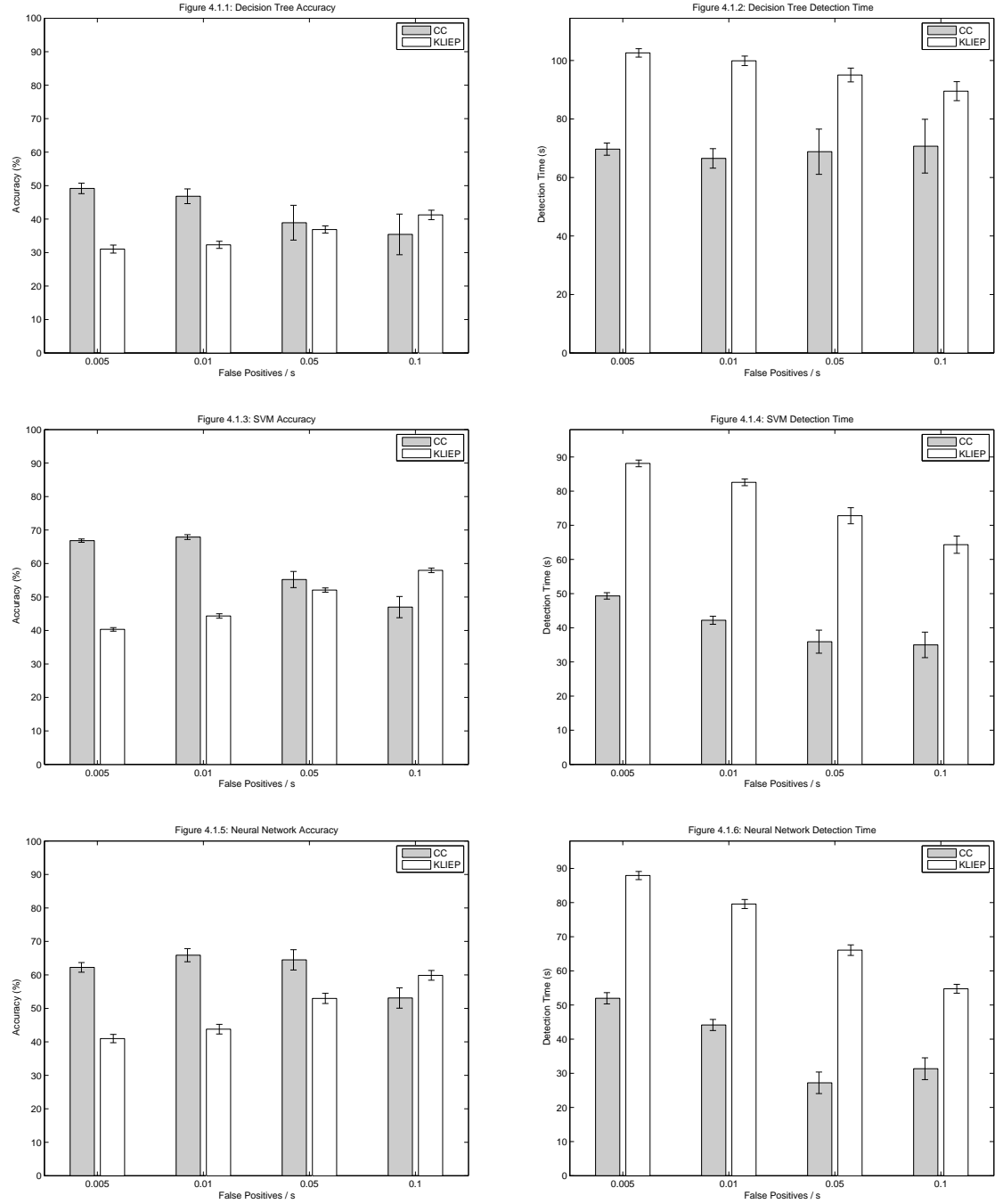


Figure 4.1: OSU Hip Results. Graphs are organized into rows by base classifier, and columns by evaluation metric. Change-point detection results were averaged over 30 splits into training, testing, and validation datasets, along with error bars showing a 95% confidence interval.

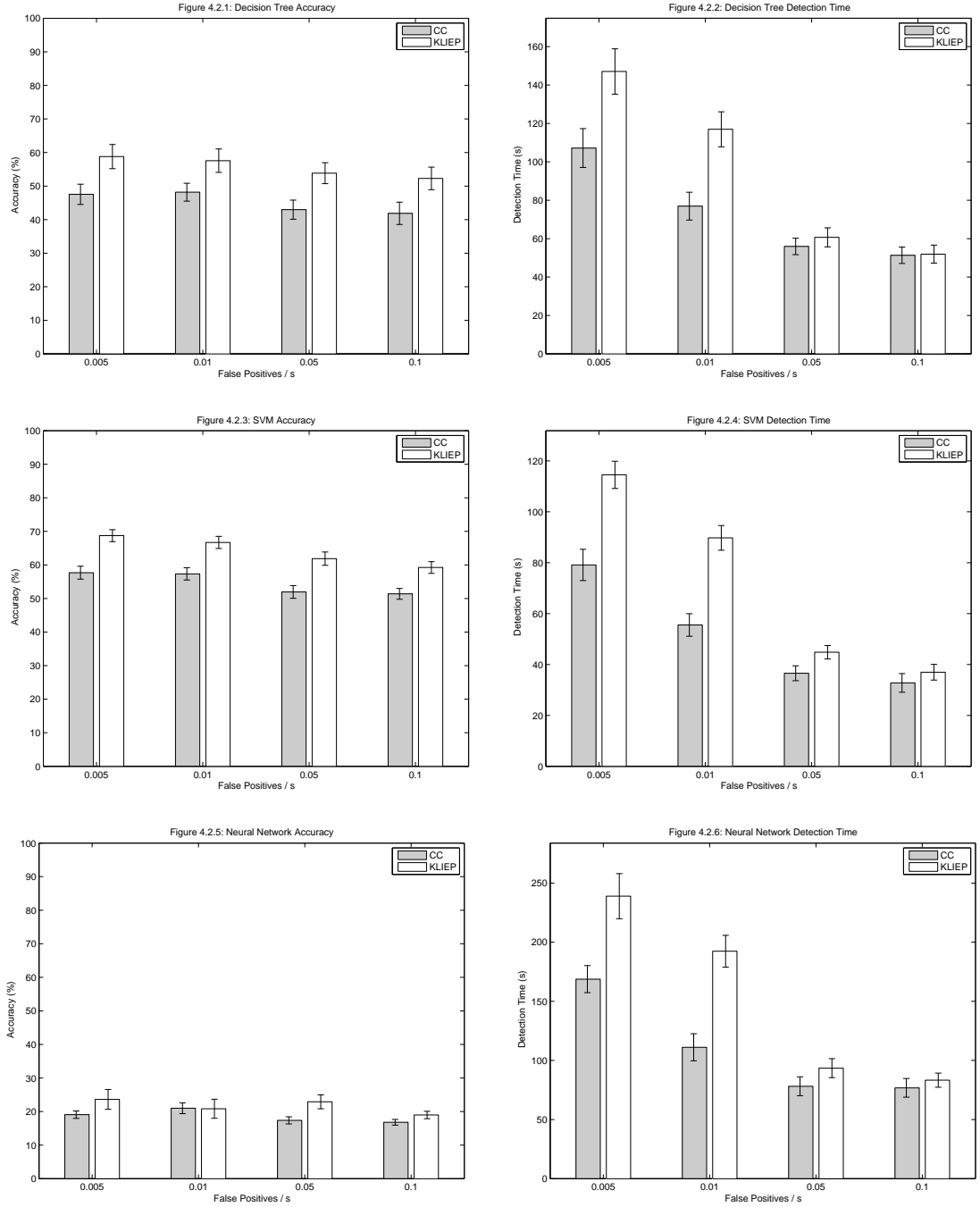


Figure 4.2: LiME Day 1 Results. Graphs are organized into rows by base classifier, and columns by evaluation metric. Change-point detection results were averaged over 30 splits into training, testing, and validation datasets, along with error bars showing a 95% confidence interval.

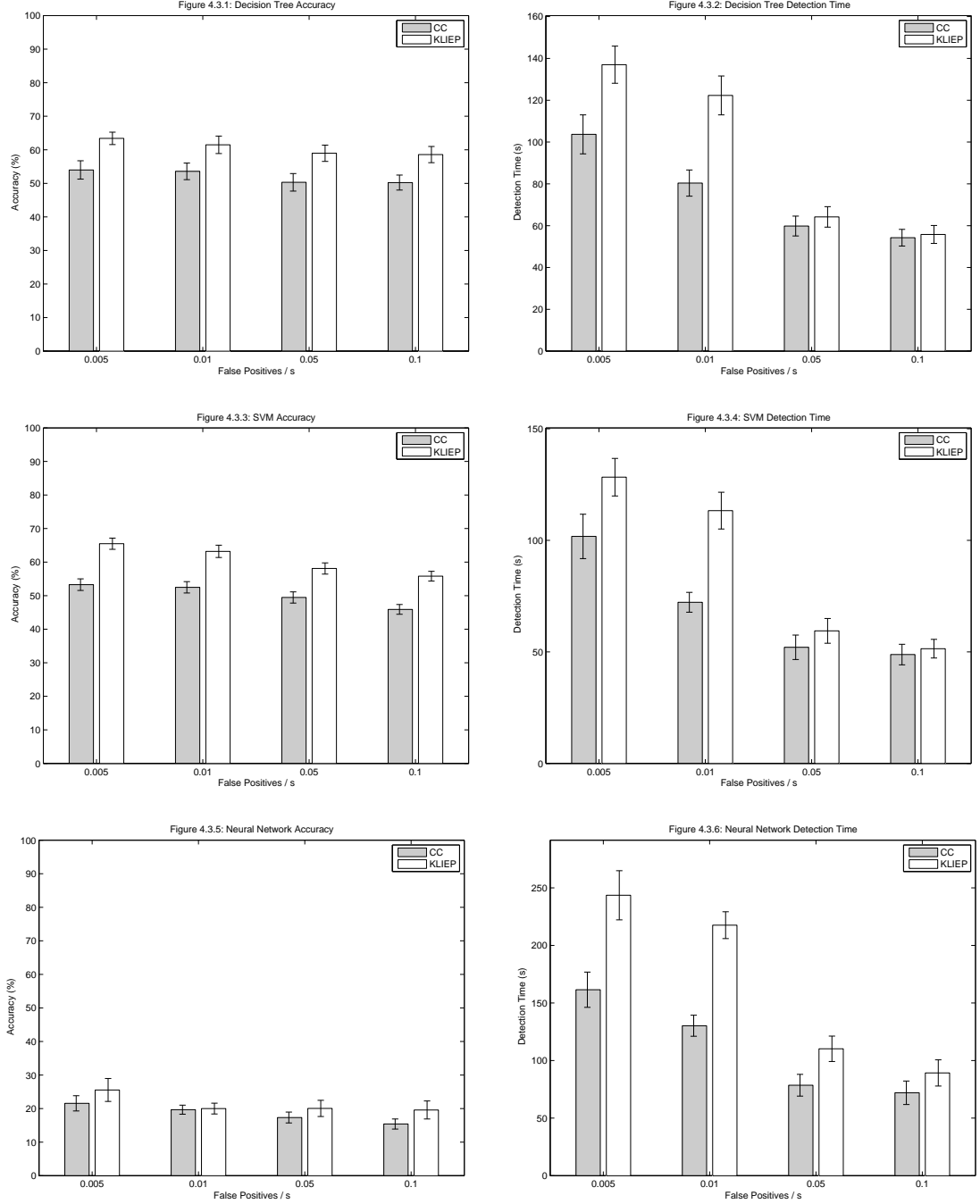


Figure 4.3: LiME Day 2 Results. Graphs are organized into rows by base classifier, and columns by evaluation metric. Change-point detection results were averaged over 30 splits into training, testing, and validation datasets, along with error bars showing a 95% confidence interval.

4.2 HMM

Results for our HMM experiments are given in Figures 4.4-4.6. Each HMM experiment was performed by splitting each time series into windows of fixed length corresponding to discrete time “ticks” in an HMM, and results for windows of length $\{10, 12, 14, 16, 18, 20\}$ seconds are shown.

For both the SVM and decision tree classifiers, accuracy was high and detection time was low across all three datasets. Accuracy and detection time were also stable with respect to window size. Further experiments on the OSU Hip dataset showed that the HMM when paired with these classifiers tends to be stable with window sizes that are greater than a few seconds, which seems to be the amount of time required to be informative. Neural networks performed somewhat more poorly and erratically across the board.

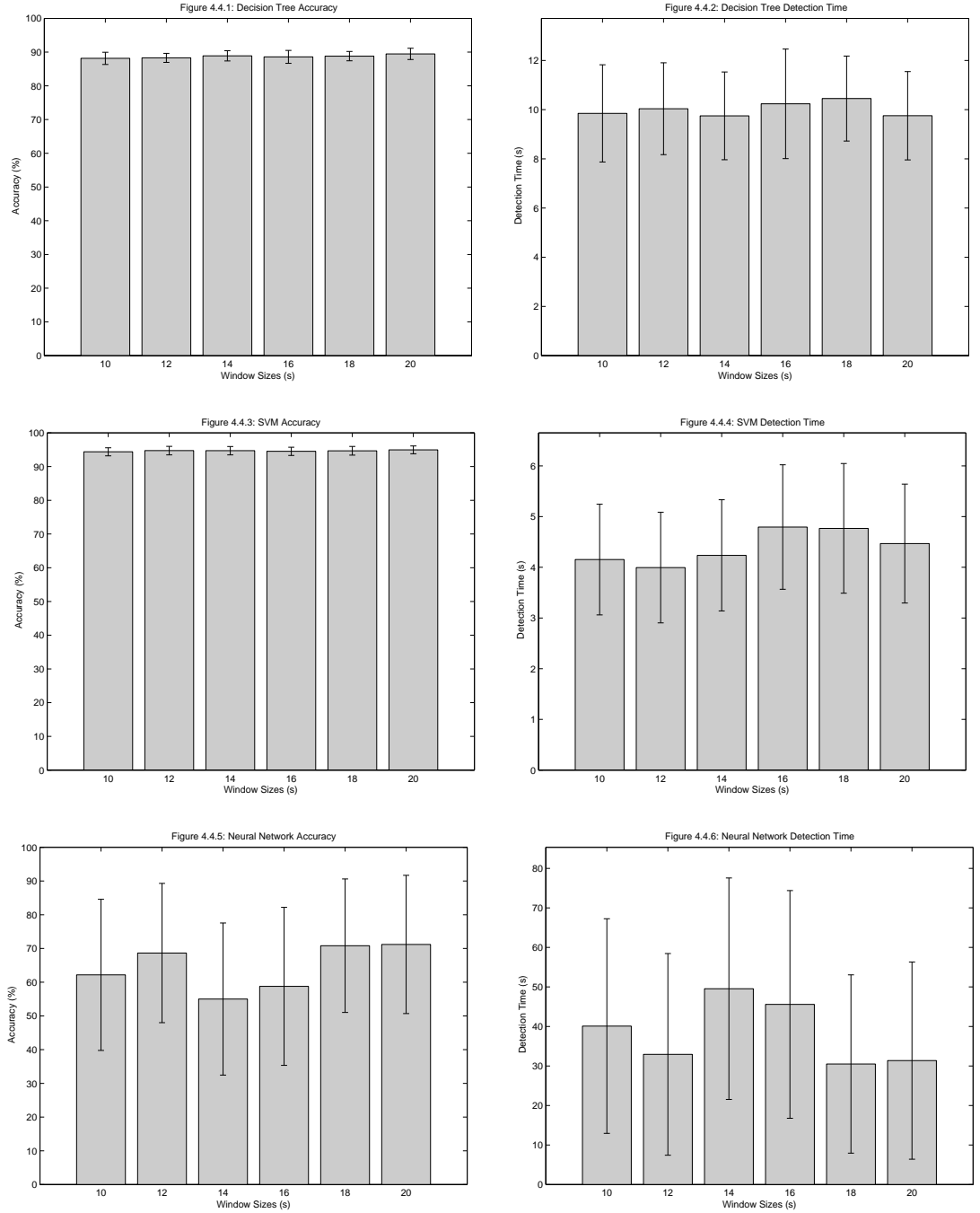


Figure 4.4: OSU Hip HMM Results. Graphs are organized into rows by base classifier, and columns by evaluation metric. HMM results were averaged over 10 splits into training (base classifier), validation, training (HMM), and testing datasets, along with error bars showing a 95% confidence interval.

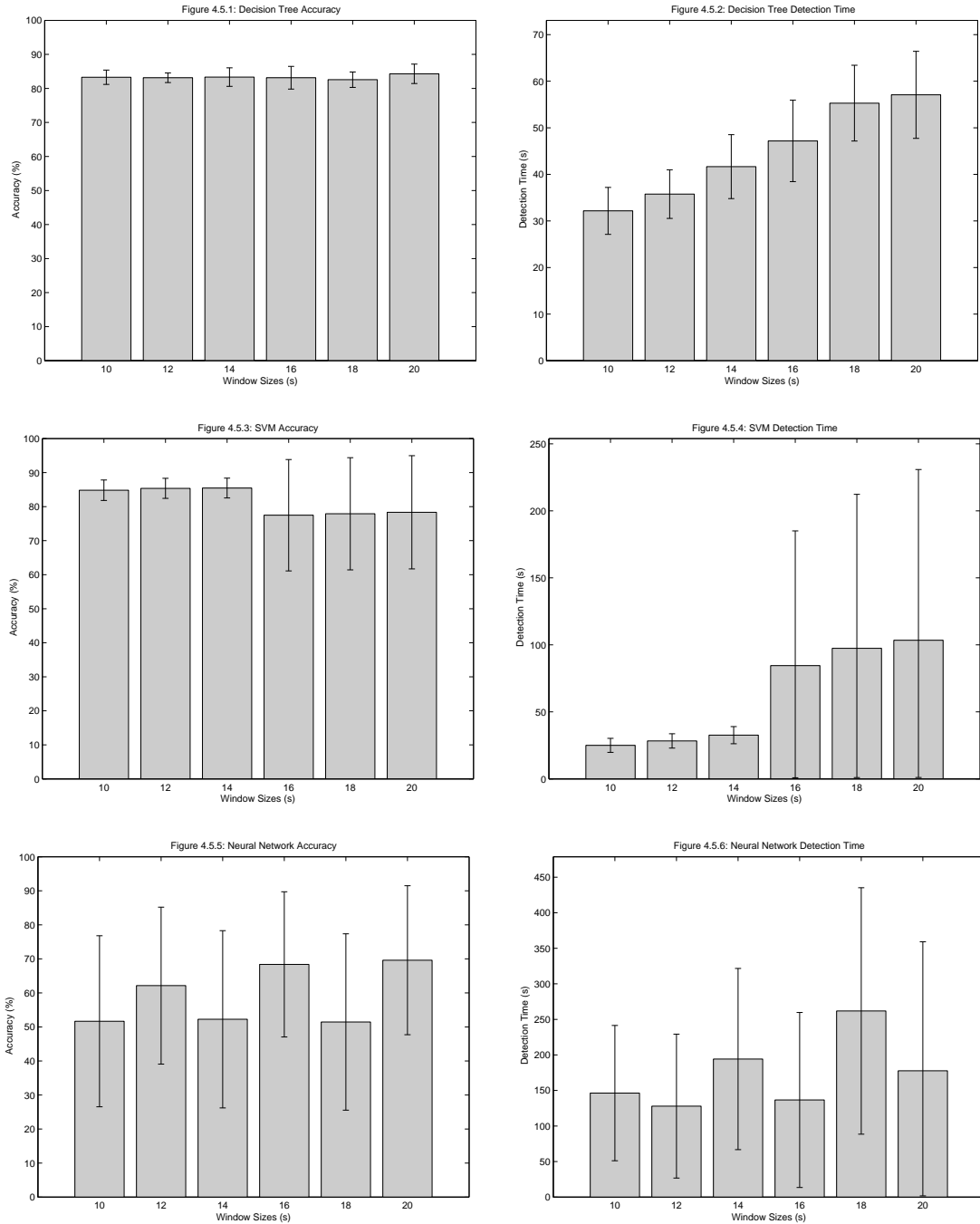


Figure 4.5: LiME Day 1 HMM Results. Graphs are organized into rows by base classifier, and columns by evaluation metric. HMM results were averaged over 10 splits into training (base classifier), validation, training (HMM), and testing datasets, along with error bars showing a 95% confidence interval.

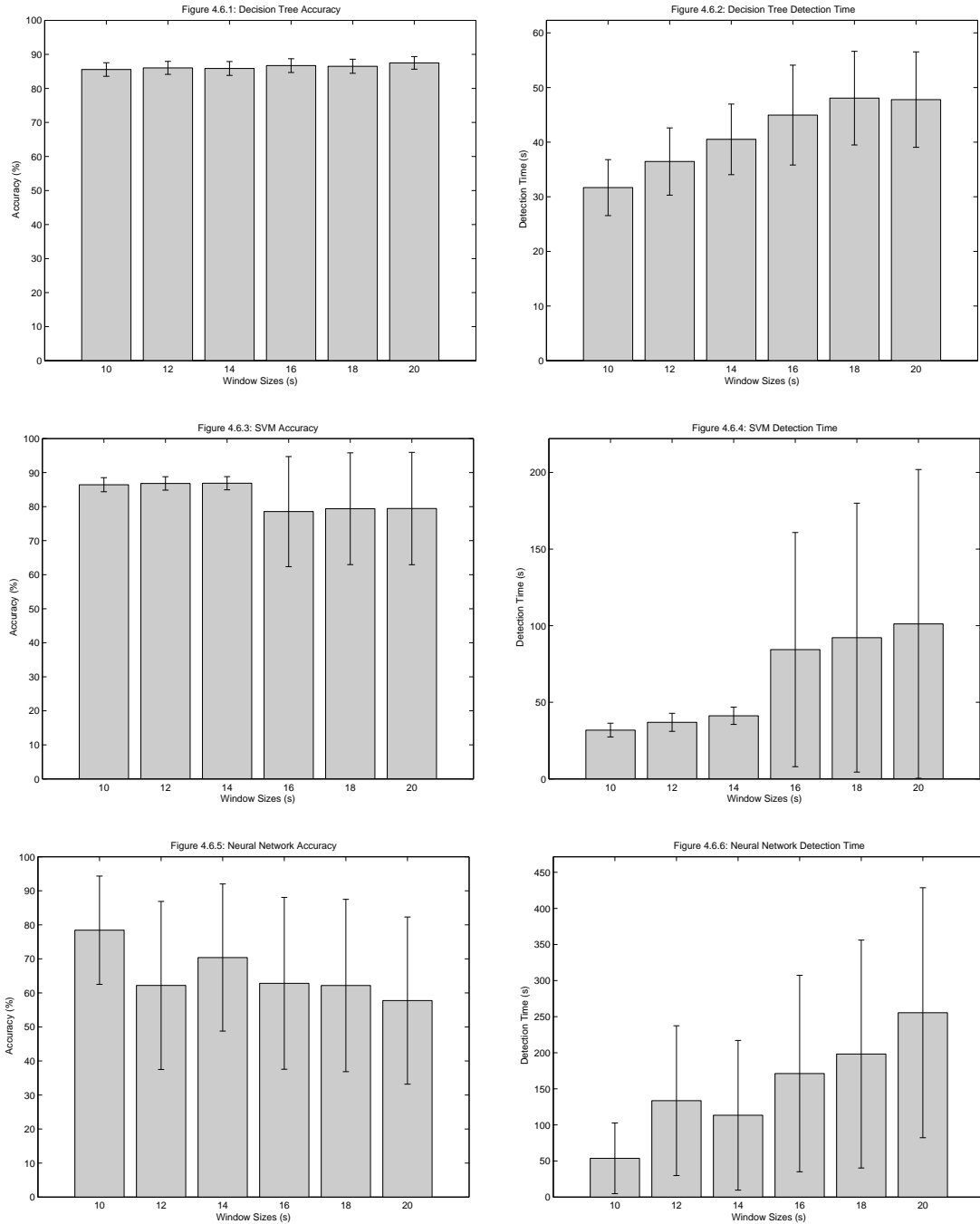


Figure 4.6: LiME Day 2 HMM Results Graphs are organized into rows by base classifier, and columns by evaluation metric. HMM results were averaged over 10 splits into training (base classifier), validation, training (HMM), and testing datasets, along with error bars showing a 95% confidence interval.

4.3 Discussion

Our results clearly show that the HMM approach outperformed the change-point detection approach, both in terms of accuracy and detection time, regardless of the dataset and base classifier. HMMs likely did better because they are trained on labelled data. By contrast, change-point detection algorithms are unaware of true activity labels—they merely signal a change when they detect a distance or dissimilarity between reference and test data. As such, change-point detection algorithms have trouble accurately segmenting free-living data, and noisy segmentation results in poor performance.

A contributing factor to the particularly high accuracy and low detection time results attained for the OSU Hip experiments was that the data consisted of activities that were synthetically glued together. The same group of activities were performed in the same order by each of the 50 subjects in this dataset, making transitions from one activity to the other very predictable for a temporal model. By contrast, the LiME datasets consisted of unsynthetic data gathered from a large set of unstructured and variable-length activities, so the activity transitions were not as predictable and are more indicative of an application of our techniques in the real world.

A final point of interest was that SVM clearly outperformed the other two base classifiers, and that the faster and simpler decision tree model did fairly well against neural networks. This result is significant because much of the previous research that has formulated activity detection as a supervised learning problem has used neural networks exclusively.

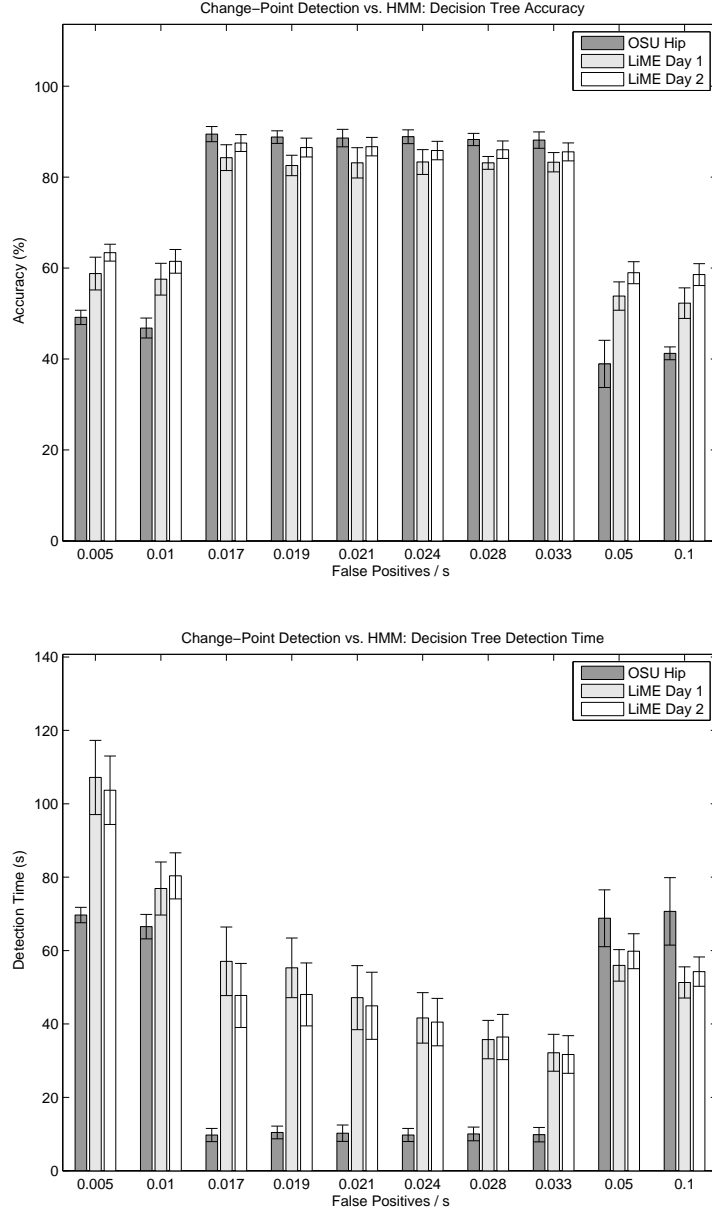


Figure 4.7: Comparison of change-point detection and HMMs in terms of accuracy and detection time. For a given false positive rate, change-point detection algorithms split a time series into windows of a certain average size, which decreases as the false positive rate increases. From this it is possible to relate false positive rates from the CPD experiments to the window sizes of the HMM experiments. Here is the HMM data with false positive rates of $\{0.033, 0.028, 0.024, 0.021, 0.019, 0.017\}$ that would corresponds to average window sizes of $\{10, 12, 14, 16, 18, 20\}$. Both graphs show decision tree results, and both show results from each of the three datasets clustered together over a single false positive rate.

Chapter 5: Conclusion

The purpose of this work was to test the feasibility of using change-point detection techniques for deciding when one activity ended within a time series and the next began, and to contrast this technique with an HMM approach. The bottom-up approach clearly outperformed the top-down approach. We also showed that the performance of a change-point detection algorithm was highly dependent on how well the data fit the modeling assumption of the algorithm, so it is plausible that a change-point detection algorithm with a modeling assumption that is in accord with the given data will perform well.

Bibliography

- [1] Gregory D. Abowd, Anind K. Dey, Peter J. Brown, Nigel Davies, Mark Smith, and Pete Steggles. Towards a better understanding of context and context-awareness. In *Proceedings of the 1st international symposium on Handheld and Ubiquitous Computing*, HUC '99, pages 304–307, London, UK, UK, 1999. Springer-Verlag.
- [2] K. Aminian, P. Robert, E. Jequier, and Y. Schutz. Estimation of speed and incline of walking using neural network. *Instrumentation and Measurement, IEEE Transactions on*, 44(3):743–746, 1995.
- [3] Suk Joo Bae, Byeong Min Mun, and Kyung Yong Kim. Change-point detection in failure intensity: A case study with repairable artillery systems. *Computers and Industrial Engineering*, 64:11–18, January 2013.
- [4] Ling Bao and Stephen S. Intille. Activity recognition from user-annotated acceleration data. pages 1–17. Springer, 2004.
- [5] S.I. de Vries, F.G. Garre, L.H. Engbers, V.H. Hildebrandt, and S. van Buuren. Evaluation of neural networks to identify types of activity using accelerometers. *Medicine & Science in Sports & Exercise*, 43(1):101, 2011.
- [6] T.V. Duong, H.H. Bui, D.Q. Phung, and S. Venkatesh. Activity recognition and abnormality detection with the switching hidden semi-markov model. In *IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 838–845, June 2005.
- [7] James Fogarty, Carolyn Au, and Scott E. Hudson. Sensing from the basement: a feasibility study of unobtrusive and low-cost home activity recognition. In *Proceedings of the 19th annual ACM symposium on User interface software and technology*, UIST '06, pages 91–100, New York, NY, USA, 2006. ACM.
- [8] Yoshinobu Kawahara and Masashi Sugiyama. Change-point detection in time-series data by direct density-ratio estimation. *Proceedings of the SIAM International Conference on Data Mining*, pages 389–300, 2009.
- [9] Jonathan Lester, Tanzeem Choudhury, Nicky Kern, Gaetano Borriello, and Blake Hannaford. A hybrid discriminative/generative approach for modeling human activities. In *In Proc. of the International Joint Conference on Artificial Intelligence (IJCAI)*, pages 766–772, 2005.

- [10] Z. Li. Exercises intensity estimation based on the physical activities healthcare system. In *Communications and Mobile Computing, 2009. CMC'09. WRI International Conference on*, volume 3, pages 132–136. IEEE, 2009.
- [11] David S. Matteson and Nicholas A. James. A nonparametric approach for multiple change point analysis of multivariate data. 2012.
- [12] David Meyer, Evgenia Dimitriadou, Kurt Hornik, Andreas Weingessel, and Friedrich Leisch. *e1071: Misc Functions of the Department of Statistics (e1071)*, TU Wien, 2012. R package version 1.6-1.
- [13] V. Moskvina and A. A. Zhigjavsky. An algorithm based on singular-spectrum analysis for change-point detection. *Communication in Statistics. Statistics and Simulations*, 32:319–352, 2003.
- [14] D.M. Pober, J. Staudenmayer, C. Raphael, and P.S. Freedson. Development of novel techniques to classify physical activity mode using accelerometers. *Medicine & Science in Sports & Exercise*, 38:1626, 2006.
- [15] Nishkam Ravi, Nikhil D, Preetham Mysore, and Michael L. Littman. Activity recognition from accelerometer data. In *In Proceedings of the Seventeenth Conference on Innovative Applications of Artificial Intelligence (IAAI)*, pages 1541–1546. AAAI Press, 2005.
- [16] Brian Ripley. *Feed-forward Neural Networks and Multinomial Log-Linear Models*, 2013. R package version 7.3-6.
- [17] M.P. Rothney, M. Neumann, A. Béziat, and K.Y. Chen. An artificial neural network model of energy expenditure using nonintegrated acceleration signals. *Journal of Applied Physiology*, 103(4):1419–1427, 2007.
- [18] Jim Rowan and Elizabeth D. Mynatt. Digital family portrait field trial: Support for aging in place. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, CHI '05, pages 521–530, New York, NY, USA, 2005. ACM.
- [19] Stewart Russell and Peter Norvig. *Artificial Intelligence: A Modern Approach*, pages 566–583. Prentice Hall, third edition, 2010.
- [20] W. Shewhart. Quality control charts. *Bell System Technical Journal*, pages 593–603, 1926.
- [21] Yoonseon Song, Seungchul Shin, Seunghwan Kim, Doheon Lee, and Kwang Lee. Speed estimation from a tri-axial accelerometer. In *Proceedings of the 29th Annual International Conference of IEEE EMBS*, August 2007.

- [22] M. Staudacher and all. A new method for change-point detection developed for on-line analysis of the heart beat variability during sleep. *Statistical Mechanics and its Applications*, 349:582–596, April 2005.
- [23] J. Staudenmeyer, D. Pober, S. Crouter, D. Bassett, and P. Freedson. An artificial neural network to estimate physical activity energy expenditure and identify physical activity type from an accelerometer. *Journal of Applied Physiology*, pages 1300–1307, 2009.
- [24] Masashi Sugiyama and all. Direct importance estimation with model selection and its application to covariate shift adaptation. *Annals of the Institute of Statistical Mathematics*, 60:699–746, 2008.
- [25] A.G Tartakovsky and all. A novel approach to detection of intrusions in computer networks via adaptive sequential and batch-sequential change-point detection methods. *IEEE Transactions on Signal Processing*, 54:3372–3382, September 2006.
- [26] G. Thatte and all. Parametric methods for anomaly detection in aggregate traffic. *IEEE/ACM Transactions on Networking*, 19:512–519, April 2011.
- [27] Terry Therneau, Beth Atkinson, and Brian Ripley. *Recursive Partitioning*, 2013. R package version 4.1-1.
- [28] S.G. Trost, W.K. Wong, K.A. Pfeiffer, and Y. Zheng. Artificial neural networks to predict activity type and energy expenditure in youth. *Medicine and Science in Sports and Exercise*, pages 1801–1809, September 2012.
- [29] Tim van Kasteren, Athanasios Noulas, Gwenn Englebienne, and Ben Kröse. Accurate activity recognition in a home setting. In *Proceedings of the 10th international conference on Ubiquitous computing*, UbiComp ’08, pages 1–9, New York, NY, USA, 2008. ACM.
- [30] Junbo Wang, Zixue Cheng, Mengqiao Zhang, Yinghui Zhou, and Lei Jing. Design of a situation-aware system for abnormal activity detection of elderly people. In *Proceedings of the 8th international conference on Active Media Technology*, AMT’12, pages 561–571, Berlin, Heidelberg, 2012. Springer-Verlag.
- [31] Y. Zheng, W.K. Wong, X. Guan, and S. Trost. Physical activity recognition from accelerometer data using a multi-scale ensemble method. In *Twenty-Fifth Annual Conference on Innovative Applications of Artificial Intelligence*. IAAA, 2013.
- [32] Yonglei Zheng. Predicting activity type from accelerometer data. Master’s thesis, Oregon State University, August 2012.

