
Undergraduate Honors Theses

Theses, Dissertations, & Master Projects

4-2017

Saving Babies Using Big Data

Evan Dienstman

College of William and Mary

Follow this and additional works at: <https://scholarworks.wm.edu/honorsthesis>



Part of the [Applied Statistics Commons](#)

Recommended Citation

Dienstman, Evan, "Saving Babies Using Big Data" (2017). *Undergraduate Honors Theses*. Paper 1048.
<https://scholarworks.wm.edu/honorsthesis/1048>

This Honors Thesis is brought to you for free and open access by the Theses, Dissertations, & Master Projects at W&M ScholarWorks. It has been accepted for inclusion in Undergraduate Honors Theses by an authorized administrator of W&M ScholarWorks. For more information, please contact scholarworks@wm.edu.

Saving Babies Using Big Data

A thesis submitted in partial fulfillment of the requirement
for the degree of Bachelor of Science in Mathematics from
The College of William and Mary

by

Evan Dienstman

Accepted for _____
(Honors, High Honors, Highest Honors)

Prof. John Delos, Director

Prof. Leah Shaw

Prof. Junping Shi

Prof. Daniel Vasiliu

Williamsburg, VA
April 21, 2017

William & Mary Mathematics Honors Project



Saving Babies Using Big Data

Student

Evan Dienstman

eddienstman@email.wm.edu

Academic Advisers

John Delos, jbdelo@wm.edu

Daniel Vasiliu, dvasiliu@wm.edu

May 10, 2017

Contents

List of Figures	5
Abstract	6
Acknowledgments	7
1 Introduction	8
1.1 Medical Background	8
1.2 Prior Work	9
1.3 Data Collection	9
1.4 Notes on the Matlab Code	10
2 Heart Rate Characteristics	12
2.1 Mean	12
2.2 Variance	12
2.3 Asymmetry	12
2.4 Sample Entropy	13
2.5 Decelerations	13
2.6 HRC Subcategories	14
3 Data Organization and Preprocessing	16
3.1 HRB Files	16
3.2 Excel Files	16
3.3 Dienstman Result Files	17
3.4 CSV Files	18
4 Time Series	21
4.1 Time Series Figures	21
4.2 Time Series Discussion	22
5 Univariate Probability Density Functions	24
5.1 Kernel Density Estimation Definition	24
5.2 Univariate PDF Figures	25
5.3 Univariate PDF Discussion	26
6 Univariate Risks	28
6.1 Risk Definition	28
6.2 Univariate Risk Figures	29
6.3 Univariate Risk Discussion	29

7	Bivariate Probability Density Functions	31
7.1	Bivariate PDF Figures	31
7.2	Bivariate PDF Discussion	32
8	Bivariate Risks	33
8.1	Bivariate Risk Figures	33
8.2	Bivariate Risk Discussion	34
9	Single Variable Logistic Regression	35
9.1	Single Variable Logistic Regression Definition	35
9.2	Single Variable Logistic Regression Figures	36
9.3	Single Variable Logistic Regression Discussion	36
10	HeRO Score	38
10.1	HeRO Score Definition	38
10.2	HeRO Score Figures	39
10.3	HeRO Score Discussion	40
11	Conclusion	42
11.1	Discussion	42
11.2	Future Work	42
A	Additional Figures	44
A.1	Time Series Figures	44
A.2	Univariate PDF Figures	46
A.3	Univariate Risk Figures	49
A.4	Bivariate PDF Figures	50
A.5	Bivariate Risk Figures	51
A.6	Single Variable Logistic Figures	52
A.7	HeRO Score Figures	54
B	Matlab Programs	56
B.1	Result Files	56
B.2	CSV Files	62
B.3	Time Series Figures	67
B.4	Univariate PDF Figures	73
B.5	Univariate Risk Figures	76
B.6	Bivariate PDF Figures	82
B.7	Bivariate Risk Figures	86
B.8	Single Variable Logistic Figures	89
B.9	HeRO Score Figures	92
	Bibliography	98

List of Figures

1.1	RR Interval	9
2.1	Deceleration	14
3.1	RCTEvents2.xls	17
3.2	Dienstman Result File Example	18
3.3	CSV File Example	19
4.1	Individual Time Series Figure Example	22
4.2	Average Time Series Figure	22
5.1	Comparison of Histogram and Corresponding KDE	25
5.2	Univariate PDF Example	26
6.1	Univariate Risk Example	29
7.1	Bivariate PDF Example	32
8.1	Bivariate Risk Example	33
9.1	Single Variable Logistic Regression Example	36
10.1	HeRO Score Figure Example	39
10.2	Average HeRO Score Figure	40
A.1	Time Series 1	44
A.2	Time Series 2	45
A.3	Time Series 3	45
A.4	Univariate PDF 1	46
A.5	Univariate PDF 2	46
A.6	Univariate PDF 3	47
A.7	Univariate PDF 4	47
A.8	Univariate PDF 5	48
A.9	Univariate PDF 6	48
A.10	Univariate Risk 1	49
A.11	Univariate Risk 2	49
A.12	Bivariate PDF 1	50
A.13	Bivariate PDF 2	50
A.14	Bivariate Risk 1	51
A.15	Bivariate Risk 2	51
A.16	Single Variable Logistic Figure 1	52
A.17	Single Variable Logistic Figure 2	52

A.18 Single Variable Logistic Figure 3	53
A.19 Single Variable Logistic Figure 4	53
A.20 HeRO Score Figure 1	54
A.21 HeRO Score Figure 2	54
A.22 HeRO Score Figure 3	55

Abstract

Because of their underdeveloped immune systems, premature babies are at an increased risk to contract many illnesses. Thus, early detection of a disease is vital to saving a premature baby's life. Current methods of detecting illnesses, however, have been inadequate, providing many false positives and insufficient amount of warning time. However, patterns in the heart rate of babies have shown signs of predicting the onset of sepsis in premature infants. Research conducted by Prof. John Delos and others suggest that low variability and clusters of decelerations in an infant's heart rate may indicate an impending septic event. Additionally, there is weak evidence that low variability may be linked to gram-positive bacteria and clusters of decelerations may be linked to gram-negative bacteria. If this statement is true, then not only will the heart rate of an infant predict the onset of sepsis, but also provide a partial diagnosis and thereby indicate the preferred treatment for the baby. However, much more work needs to be done to prove this hypothesis. Over twelve terabytes of data has been collected on premature babies' heart rate and breathing. To search through this data, one first needs to know what to look for. Unfortunately, only looking for low variability and clusters of decelerations would be inadequate since most babies experience some low variability and decelerations in their heart rate at some point. Therefore, sophisticated statistical analysis is necessary to quantify this data. The general idea of this analysis includes creating many different heart rate characteristics (HRCs) and measuring their predictive power through multiple methods. The results of our research indicate that the HRCs of variance, sample entropy, and asymmetry are strong predictors of illness. However, no HRC shows strong signs of indicating the type of invading organism that caused the illness.

Acknowledgments

I would like to thank my academic advisors at William & Mary, Prof. John Delos and Prof. Daniel Vasiliu, for their instruction and support during this entire process. Various Matlab programs used in the research are original/modified programs by Prof. Douglas Lake and Abigail Flower. I would also like to thank our UVA friends, Prof. Douglas Lake and Dr. Randall Moorman, for their guidance on the medical side of this project. The data calculations in this report were performed on computational resources supported by William & Mary's high-performance computers (SciClone). This project is an extension of work started as part of my participation in EXTREEM-QED program at William & Mary. The EXTREEM-QED program is supported by NSF Grant 1331021.

Chapter 1

Introduction

Premature babies are at an increased risk of having a septic event due to their immature immune systems. Therefore, nurses must frequently take blood samples of the infants in order to determine a patient's health status. Blood work, however, is a slow process, so doctors might not be performing the best treatment for the infant while waiting for the results of a test. Furthermore, drawing blood, as with any invasive measure, could cause complications. Our goal is to try to use the baby's heart rate as a quick, noninvasive way to decide if an infant is unhealthy. Research has already shown some signs that reduced variability and transient decelerations in the heart rate could predict sepsis [4]. We also hypothesize that low variability may be linked to gram-positive bacteria and clusters of decelerations may be linked to gram-negative bacteria, giving insight into the preferred treatment for the baby. In order to test our two hypotheses, we will develop heart rate characteristics (HRCs) to better understand the data. We will then test the usefulness of the HRCs at predicting sepsis. In the end, the ultimate goal will be to test if our HRCs are predictive of sepsis and the hypothesis that low variability is associated with gram-positive diseases and clusters of decelerations are associated with gram-negative diseases.

1.1 Medical Background

1.1.1 Sepsis

Sepsis refers to an inflammatory response caused by the body fighting off an infection [8]. In order to detect sepsis, a nurse will look for changes in a baby's body temperature, digestion, and other vital signs. Unfortunately, these symptoms may not appear until long after the baby has become infected. Another challenge is that even though early treatment with antibiotics is very effective in treating sepsis, without the proper diagnosis of the invading organism, a doctor will not be able to administer the best drug to the patient. A nurse, consequently, must take a blood sample in order to properly identify the pathogen. However, this procedure is time consuming and not ideal for very small infants. In our report, we will focus on five classifications of invading organisms: coagulase-negative staphylococci (CONS), gram-positive bacteria, gram-negative bacteria, fungus, and other. The hope from the heart rate analysis is two-fold: 1) give earlier warning than current methods about a septic event and 2) provide information about the type of organism that caused the septic event. If we accomplish both our goals, then not only will we be able to provide doctors with an inexpensive, noninvasive way of detecting sepsis in its early stages, but also provide a partial diagnosis, and consequently a treatment, for the pathogen.

1.1.2 RR Intervals

An electrocardiogram (EKG) displays a time series of electric impulses created by the heart. In this time series, the largest peaks occur at the beginning ventricular contraction. This peak is defined as the R-peak. Physicians use the R-peak to represent the time a heart beat occurs. Therefore, the RR interval is the time between one R-peak to the next, or rather, the time between heart beats. When analyzing heart rate throughout this report, we will use the RR intervals. Note that even though the the RR intervals are not technically a rate, one can easily obtain a heart rate by finding the number of intervals in a certain unit of time. Another way to relate the RR intervals to a heart rate is that large RR intervals imply a slower heart rate and small RR intervals imply a faster heart rate. [1] Figure 1.1 shows an example time-series of an EKG with one RR interval.

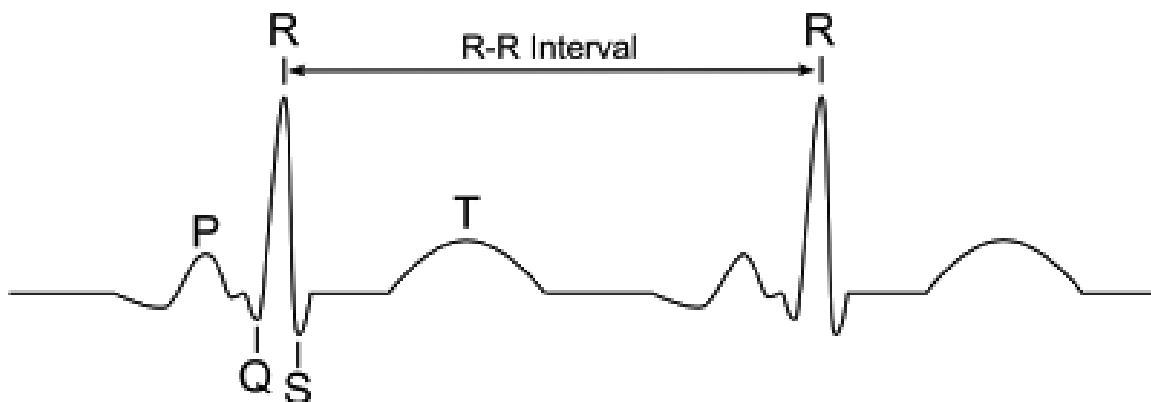


Figure 1.1: RR Interval[9]

1.2 Prior Work

As reported in *Mortality Reduction by Heart Rate Characteristic Monitoring in Very Low Birth Weight Neonates: A Randomized Trial* by Dr. J. Randall Moorman and Prof. Douglas E. Lake, et al., four heart rate characteristics were found to provide early warnings of illness: variance, asymmetry 1, asymmetry 2, and sample entropy. We will explain the definitions of these measures in Chapter 2. A Heart Rate Observation score, or HeRO score, which was previously derived using a sample of a few hundred sepsis events, was also used in the randomized clinical trial to provide warning of illness. The HeRO score is a measure that combines the four heart rate characteristics. We will explain more about the HeRO score in Chapter 10. Specifically, the trial showed that as a consequence of these early warnings, mortality in neonatal intensive care units was reduced by between 10 – 40% amongst various cohorts.

In this paper, we reanalyze the data from the clinical trial and develop a new HeRO score based off this data. We also incorporate additional measures of heart rate variability into the score. We hope that these two changes will improve the statistical predications of the HeRO score.

1.3 Data Collection

All data was collected from electrocardiograms at eight neonatal intensive care units

(NICUs) during the randomized trial led by Dr. J. Randall Moorman [6]. The locations of the NICUs have been withheld to protect the identity of the patients. In total, the trial included around 3,000 babies with each baby having roughly one month’s worth of data. From the 3,000 babies, about 1,000 septic events were recorded. Note that it is possible for one baby to experience multiple septic events.

Each infant has a file containing a vector of RR intervals as well as the time each RR interval ends. Times are calculated from time 0 corresponding to when the patient was first connected to the EKG. Because the RR intervals are times themselves, the time vector is simply a cumulative sum of the RR interval vector assuming the monitor is never disconnected. However, the time vector is necessary since babies may be disconnected from the monitor at times. We, therefore, need the time vector to indicate at what time each RR interval occurred in the patient’s history. An example of these two vectors are shown below:

```
rr_intervals = [420.25, 421.75, 420.25, 418.25, 422.00, 423.75 ...]
tt_intervals = [420.25, 842.00, 1262.3, 370470, 370880, ...]
```

In the example above, `rr_intervals` refers to the RR intervals and `tt_intervals` refers to the times of the RR peaks. Note that `tt_intervals` is a cumulative sum of `rr_intervals` until the large jump from 1,262.3 to 370,470. A large jump like this corresponds to a time when a baby was disconnected from the monitor. Smaller jumps may also occur when the analyzing software fails to find a beat. In this example, both `rr_intervals` and `tt_intervals` are measured in milliseconds. Throughout the rest of the report, we measure the RR intervals in milliseconds and the times of the RR intervals in days.

The RR intervals and the times of the RR intervals constitute the bulk of the data for a patient. Other data for a patient includes a site number corresponding to the patient’s NICU, an ID number, demographic information, and a file indicating the times when the patient was ventilated. If a patient is having difficulty breathing, the nurses will ventilate the infant. We keep a record of this information because we believe ventilation reduces heart rate variability. Lastly, we have another file containing all the patients that had a septic event, the time of the event, and the invading organism that caused the event. Note that a single patient can experience multiple septic events while at a NICU.

In order to better understand the data, we break up all the `rr_interval` vectors into half hour pieces. Later, when we calculate statistical measures on the RR intervals, we calculate them for each half hour vector instead of over the entire vector. Doing our calculations this way helps us determine change in the various measurements over time. The choice to divide our vectors into half hours is an empirical decision aided by opinions of Prof. Lake and Dr. Moorman at UVA. Note that changing the length of the vector pieces will have an effect on the statistical measurements, and similar work in the field uses 5 or 10 minute intervals. However, for the purpose of this report, we use 30 minute intervals.

1.4 Notes on the Matlab Code

Throughout this thesis, we will reference Matlab code used to create data files and figures. To help relate the code to the report, each section using code will contain the names of the files involved at the beginning. Consider the following example below:

```
Bivariate Risk Figures: multiple_bivariate_risk_figures →
one_bivariate_risk_figure → one_bivariate_risk_plot → one_risk_matrix →
one_prob_matrix
```

The text above shows the code needed to create the bivariate risk figures. Note that

`multiple_bivariate_risk_figures` → `one_bivariate_risk_figure`

means that `multiple_bivariate_risk_figures` calls `one_bivariate_risk_figure`. However, sometimes an arrow does not mean one function calls another, but rather the second function is the next function in the procedural order. While the report reviews important aspects of the code, all code will not be discussed in detail. However, Appendix B contains all the Matlab code for reference written by Evan Dienstman. Code written by other people has been omitted from this report.

Chapter 2

Heart Rate Characteristics

Now that we have RR intervals broken up into half hour sections, we need to develop statistical measures for each half hour that may be predictive of illness. We will focus on seven statistical measures, or heart rate characteristics (HRCs), throughout this report: mean, variance, sample entropy, three measures of asymmetry, and decelerations. Additionally, we will calculate five subcategories for each HRC: the raw HRC value, the 10th percentile of the HRC over the past 12 hours, the 50th percentile of the HRC over the past 12 hours, the 90th percentile of the HRC over the past 12 hours, and the slope of the HRC over the past 2 days. Accordingly, with seven HRCs and five subcategories for each one, we get a total of 35 HRCs we will analyze. The remainder of the chapter gives the definitions for each of these 35 HRCs.

2.1 Mean

The first HRC is the mean. For each half hour, we simply calculate the average RR interval length. The mean for each half hour typically ranges from 420 *ms* - 440 *ms*.

2.2 Variance

Like the mean, we also calculate the variance for each half hour. However, when doing our calculation, we take the natural logarithm of the variance so the scale is easier to visualize. These log-variance values range from 0 - 5.

2.3 Asymmetry

Asymmetry is a measure of how skewed the data looks. To calculate asymmetry, consider a vector with N RR intervals with median m where B is the set of intervals below m and A is the set of intervals above m . We can then define the following quantities:

$$r_1 = \frac{1}{N} \sum_{i \in B} (RR_i - m)^2 \quad (2.1)$$

$$r_2 = \frac{1}{N} \sum_{i \in A} (RR_i - m)^2 \quad (2.2)$$

Notice that the calculations for r_1 and r_2 are similar to the calculation of variance except with the median instead of the mean. Finally, we define three measures of asymmetry:

$$asymmetry_1 = \ln(r_1) \tag{2.3}$$

$$asymmetry_2 = \ln(r_2) \tag{2.4}$$

$$asymmetry_ratio = \frac{asymmetry_2}{asymmetry_1} \tag{2.5}$$

For each half hour, typical `asymmetry_1` and `asymmetry_2` values range from 0 - 10 and `asymmetry_ratio` values range from 0 - 3. Any `asymmetry_ratio` greater than 1 indicates that there are more intervals greater than the median than are less than the median, i.e., the data is skewed towards large intervals, indicating a slow heart rate [5].

2.4 Sample Entropy

The exact definition of sample entropy is beyond the scope of this report. However, the general idea is that sample entropy is a measurement of how random the numbers in the RR interval vector occur. A low sample entropy means the RR interval vector is fairly regular while a high sample entropy means the intervals appear to be random. Values of sample entropy range from 0 - 1 with 0 indicating the signal is completely periodic and 1 indicating the signal is completely random. Note that a horizontal line would be considered perfectly periodic and would have a sample entropy of 0. A sine or cosine function would also have a very low sample entropy. For a complete explanation of sample entropy, please refer to the paper *Physiological time-series analysis using approximate entropy and sample entropy* by J.S. Richman and J.R. Moorman [7].

2.5 Decelerations

The final raw HRC is number of decelerations. We define a deceleration as a sharp increase in the RR intervals followed by a sudden return to a baseline. Recall that an RR interval is the time between heart beats; therefore, an increase in the RR intervals represents a decrease in the heart rate. Figure 2.1 shows an example of a deceleration highlighted in red over time series of RR intervals.

Figure 2.1 contains one half hour worth of RR intervals. Within this half hour, there are many more decelerations than one highlighted in red. The computer algorithm for detecting decelerations was developed by Abigail Flower [2]. The algorithm uses a template deceleration, such as the one highlighted in red, and sweeps the template through the signal. The algorithm then makes a decision if the part of the signal in question is a significant deceleration based on the height of the peak and how well the signal matches the template. The algorithm then records the number of decelerations found. Typical values for the number of decelerations in a half hour range from 0 - 10 with rare cases going up to 100.

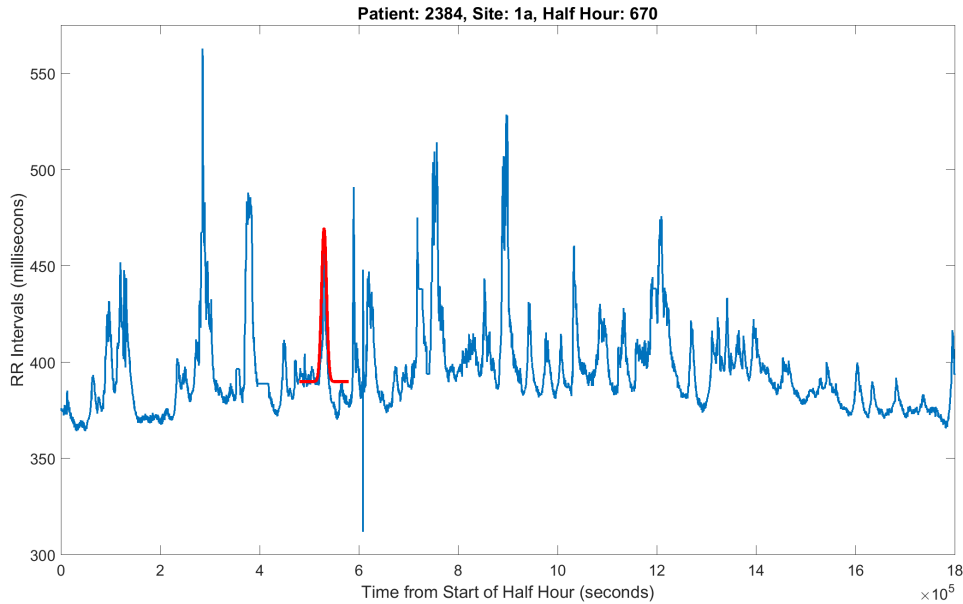


Figure 2.1: Deceleration

2.6 HRC Subcategories

In addition to all seven raw HRC measurements described above, we also calculate four more measurements for each HRC. The first three are the 10th, 50th, and 90th percentiles of the HRC values over the past 12 hours. To better understand the procedure for calculating these subcategories, consider the example of calculating the 10th percentile of variance. If all half hours are perfectly sequential in terms of their start times, *i.e.* no missing data, a window of 12 hours in the past will correspond to a vector of 24 half hours. We then calculate the 10th percentile for these 24 variance values and record this value as the 10th percentile of variance for the current half hour. The logic behind taking the 10th percentile is that one very low variance over a 12 hour window might be a stronger predictor of sepsis than looking at all the values for each half hour. However, because babies are frequently disturbed by the nurses, we take the 10th percentile instead of the minimum to avoid taking outliers caused by outside influences. Note that for a vector of 24 values, the 10th percentile is about the second lowest value.

We then repeat this process for next half hour which will result in shifting our 12 hour window up one half hour. Consequently, this new 24 half hour vector will be identical to the previous one with the exception of the oldest half hour being removed from the end and the current half hour being added to the beginning. We, therefore, should expect to see similar 10th percentile values for consecutive half hours most of the time and sudden changes sometimes once a new very low variance is captured by the window. The 50th and 90th percentiles follow a similar procedure.

The last subcategory is slope. In order to calculate the slope, we first index the raw HRC values over the past 2 days. We then take a linear fit of the points in this vector and map the slope of this fit to the HRC slope for the current half hour. Our thought process behind recording the slope is that the actual value of the HRC might not be as strong of a predictor as the rate of change of the values.

For any HRC subcategory, if there is insufficient data as a result of a baby being newly connected to a monitor or from mechanical errors, we mark the HRC as “not a number” or *NaN* in Matlab. For exact thresholds of how much data we need, please see the documentation of the relevant code.

Chapter 3

Data Organization and Preprocessing

In this section, we will focus on where we get the data for our HRC calculations and how we store the results of those calculations.

3.1 HRB Files

All the information from the EKG monitors comes from `.hrb` files. Each baby has one HRB file for his/her entire stay in the NICU. The HRB files contain the RR intervals as well as the time of each RR interval. Additionally, the monitor also reports a 0 or 1 for each interval to indicate whether the interval was good (0) or bad (1). Bad intervals may be the result of the monitor missing a heartbeat or a mechanical error. Lastly, each file is labeled with the four digit ID of the baby and contained in a folder with a site number of his/her NICU. The eight site numbers are 11, 13, 15, 23, 24, 26, 27, and 30. We will discuss what we do with this information in Section 3.3.

3.2 Excel Files

Apart from the HRB files, all the other information we have is contained in Excel files. An example of what one of these Excel files looks like is shown in Figure 3.1. This file, called `RCTEvents2.xls` contains information for every baby who had a septic event. Moving from left to right, the columns in the file are ID, site, birth weight, gestational age, group, type of organism that caused the septic event, days of age at the time of the event, death within 14 days after the event, death within 30 days after the event, HeRO score at the time of the event, and ventilation status at the time of the event. For the purpose of this report, we will ignore the group column. The organism column can contain five different classifications of organisms: 1 - coagulase-negative staphylococci (CONS), 2 - gram-positive bacteria, 3 - gram-negative bacteria, 4 - fungus, and 5 - other. The HeRO score is a type of logistic regression that we will discuss in further detail in Chapter 10. Concerning ventilation status, we have another Excel file containing specific times when each patient was ventilated. Later, when we need to determine the ventilation status of each individual half hour, we will use this file for the information. Again, we keep a record of the ventilation status because we hypothesize that babies have lower heart variability

while ventilated. Finally, one other Excel file contains demographic information about all babies that we will use later for separating babies into categories.

	A	B	C	D	E	F	G	H	I	J	K
	ID	SITE	BWT	EGA	Group	Organism	DaysofAge	Death14Days	Death30Days	HeRo	Ventilated
1	ID	SITE	BWT	EGA	Group	Organism	Days of Age	Death<14 ...	Death<30D...	HeRo	Ventilated
2	4102	11	865	28	2	1	14.8688	0	0	1.1999	0
3	4146	11	935	26	2	1	9.7931	0	0		0
4	4147	11	1223	28	1	1	19.4340	0	0		1
5	4156	11	947	27	2	1	63.4583	0	0	1.3275	0
6	4168	11	1316	33	2	1	48.3410	0	0	0.2038	1
7	4180	11	815	25	1	1	11.6368	0	0		0
8	4185	11	710	24	1	4	23.3646	0	0	1.6763	1
9	4189	11	807	25	2	2	50.0771	0	0	3.4549	1
10	4189	11	807	25	2	4	57.3104	0	0	0.9689	1
11	4202	11	800	30	2	3	61.1083	0	0	2.3973	0
12	4214	11	619	25	1	3	31.7625	0	0	1.1459	0
13	4214	11	619	25	1	1	35.9292	0	0	0.7485	1
14	4214	11	619	25	1	1	77.7799	0	0	0.4685	1
15	4215	11	540	25	2	1	7.5007	0	0	0.5952	1
16	4219	11	982	25	1	1	10.6681	0	0	0.6463	0
17	4221	11	973	26	2	3	11.8819	0	0	1.8641	0
18	4221	11	973	26	2	1	45.3493	0	0	1.3902	0
19	4221	11	973	26	2	1	74.7354	0	0		0
20	4222	11	650	24	2	1	19.1090	0	0	2.9732	1
21	4222	11	650	24	2	1	103.8292	0	0	0.2919	0
22	4254	11	810	26	2	1	8.9319	0	0	2.5591	1
23	4271	11	1100	30	2	1	10.5854	0	0	0.7148	0
24	4278	11	677	24	2	1	13.1243	0	0	4.8366	1
25	4278	11	677	24	2	3	26.8361	0	0	3.5846	1
26	4278	11	677	24	2	1	58.7563	0	0	0.2874	0
27	4294	11	824	24	1	1	89.1708	0	0	4.1154	1
28	4294	11	824	24	1	4	94.1597	0	0	0.4841	1
29	4306	11	637	24	2	4	12.0500	0	0	2.4451	1

Figure 3.1: RCTEvents2.xls

3.3 Dienstman Result Files

Dienstman Result Files: `Dienstman_submit_parallel` → `multiple_result_files` → `one_result_file`

Now that we have discussed all the raw data files, we can discuss how we process and organize the data. For each HRB file, we need to break up the information into half hours and calculate the HRCs for each half hour. The result of this procedure will be files similar to the one in Figure 3.2. The file contains a structured array with many fields. The first field is the start time of the half hour. The next 35 fields are the HRCs discussed in Section 2. The `Good_Frac` field refers to the percent of intervals in the half hour the EKG monitor marked as “good”. Because of mechanical errors or disturbances from the baby, the monitor might miss a heartbeat. Therefore, if an interval is recognized to be “bad”, an algorithm in the machine will mark it. Ideally, we want the fraction of good beats to be as close to 1 as possible. However, we currently do not remove any half hours for being below a

After calculating the HRCs for each half hour, we want to group like half hours into CSV files for easier access. Such groups include sick vs. healthy, ventilated vs. not ventilated, and the type of organism that went on to cause the illness for babies who experienced a septic event. For this report, our definition of sick is any half hour that occurs 12 hours before a recognized septic event. The definition of healthy is any half hour that is 7 days or more before an event or 3 days or more after an event. Concerning ventilated vs. not ventilated, the ventilation status of each half hour can be found within the Excel file RCT_Mechanical_Ventilation_Times.xlsx. Additionally, the file RCTEvents2.xls contains the organism that caused the illness if a baby had an event, and the file RCT_Demographics_ALL.xls contains the birth weight, gestational age, and gender of each baby.

The screenshot shows an Excel spreadsheet titled 'rhc_healthy_org_1_vent_0.csv'. The spreadsheet contains a large table with columns labeled A through Y. Column A is 'Site', B is 'ID', C is 'Half Hour', D is 'Event Time', E is 'Gender', F is 'Gestational', G is 'Good', H is 'Health', I is 'Organism', J is 'Study', K is 'Ventilated', L is 'Vent Switch', M is 'Vent Copies', N is 'Weight', and O through Y contain various 'Asymmetry' columns. The table has multiple rows of data, each representing a half-hour observation for a specific site and ID.

Figure 3.3: CSV File Example

Figure 3.3 depicts what one of these CSV files looks like. Every row in the CSV file corresponds to one half hour. Half hours are first selected from the Dienstman result files as either being sick or healthy using the file RCTEvents2.xls which contains the event times. It is important to note that the CSV files do not contain every half hour in the data set since some half hours might not be classified as either sick or healthy. The Dienstman result files contain a thorough record of this information organized by infant. Once we select a half hour, we record the site, ID, start time, and HRCs for that half hour. If the baby had a septic event, we also record the organism, event time, birth weight, and gestational age. Lastly, for both healthy and sick, we record the ventilation status for that half hour. (The last two columns in the CSV files are markers which can be ignored for the sake of this report.) In order to create the CSV files, a user calls the following commands on William & Mary’s SciClone computers:

Listing 3.2: Usage for `Dienstman_submit_batch`

```
1 >> chmod u+x Dienstman_submit_batch.txt
2 >> dos2unix Dienstman_submit_batch.txt
3 >> qsub ./Dienstman_submit_batch
```

These commands call `Dienstman_submit_batch.txt` which in turn calls `csv_files.m`. The programs make many CSV files. The first CSV file is an extensive file containing all the sick and healthy half hours. The remainder of the files only contain specific types of half hours. For example, one CSV file contains ventilated sick half hours of babies who had a septic event caused by organism 1. Breaking the CSV files down into these smaller subcategories helps with organization and loading in data when we do our analysis later. In general, the CSV files help facilitate manipulation of large data sets that would be more difficult with the result files. Thus, a majority of the statistical techniques later in the report will use the CSV files. The remaining Matlab functions in the CSV series calculate meta information from the CSV files that we will use later. For more information about how to make the CSV files, please see Appendix B.

Chapter 4

Time Series

We are finally ready to analyze the predictive power of the HRCs. Before we use the CSV files to do a more sophisticated analysis, we will first use the result files to create a time series graph of the HRCs for each infant who had a septic event. In this manner, we can gain some intuition about how the HRCs behave, especially leading up to an event for the babies that experienced one. We will also create a time series of the average HRC values relative to the time of an event. The hope is that we will see significant changes in the HRC values before an event, predicting the onset of sepsis.

4.1 Time Series Figures

Time Series Figures: `multiple_event_figures_caller` → `multiple_event_figures` →
`one_event_hrc` → `one_event_plot`

Figure 4.1 shows what a time series figure looks like for a single baby. The seven subplots represent the seven HRC. On each subplot, there are five time series for the five HRC subcategories. Note that the slope HRC uses the right y-axis. The black vertical line indicates the time of the event, and the x-axis shows the half hour index relative to the time of the event. The set of time series in a dark color depict times a baby was ventilated and light colors indicate when the baby was not ventilated. Lastly, the black horizontal line gives the average HRC value across all half hours of all babies for comparison purposes. For further examples of individual time series figures, see Appendix A where we have included examples where the baby was always unventilated, examples where the baby was always ventilated, and examples where the baby is missing data at certain times. The documentation of `multiple_event_figures_caller.m` provides more information about these figures.

After we create the time series graphs for each baby who had an event, we can then average these time series across all babies. Figure 4.2 shows the results of this averaging. To better understand the average time series, consider the data point right before the event on the variance subplot of Figure 4.2. This data point was the result of averaging all the variances at the half hour right before the event for each baby. However, since not all babies have data for this half hour, this data point was calculated using only the subset of the babies who had data at this half hour. We then repeat this process for every half hour relative to the event, noting that each average was calculated using a different amount of half hours depending on how much data we have. To get a sense about how many babies

were used for the averages, we average the number of babies used to calculate each half hour across all the average time series. This number is reported in the legend of Figure 4.2.

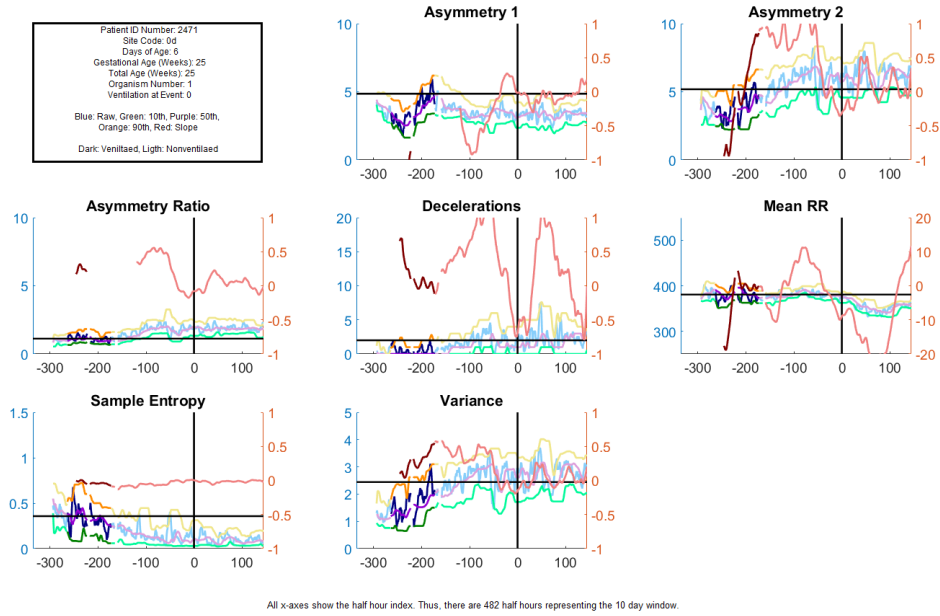


Figure 4.1: Individual Time Series Figure Example

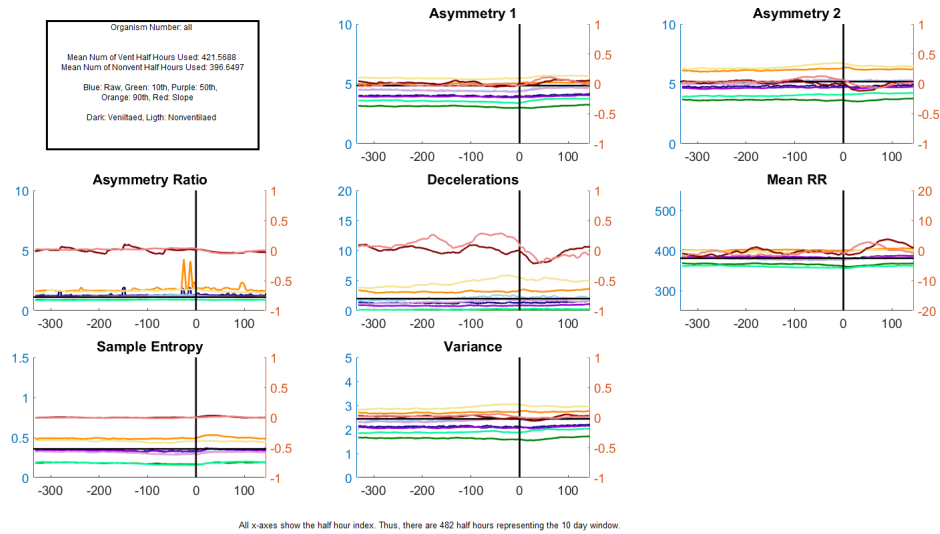


Figure 4.2: Average Time Series Figure

4.2 Time Series Discussion

The major result from Figure 4.2 is that all of the average time series are relatively flat. We notice slight changes in the HRCs before an event, but nothing that gives us a strong indication that an HRC is very predictive. The individual plots can help explain this

result. Most seem to fluctuate quite a bit, with some babies having very low HRCs and some babies having very high HRCs.

While the time series plots are useful references to visualize the data for a particular baby, they do not provide an adequate method for predicting sepsis. We therefore must resort to more sophisticated techniques, which we will discuss in the next section.

Chapter 5

Univariate Probability Density Functions

We can now use the CSV files to make probability density functions (PDFs) of the HRCs. By comparing PDFs from different groups of half hours, we can gain insight into which HRCs are more predictive. For example, recall that we have now separated the half hours into a sick category and healthy category. Using methods we will describe below, we can create a PDF for a particular HRC using only the sick half hours and a PDF for the same HRC using only the healthy half hours. If the two PDFs look significantly different, we can conclude that the HRC in question is a predictor of sepsis. We then repeat this process for every other HRC.

The naive approach for analyzing any data set is to first make a histogram. The most basic form of this approach is a histogram for one variable. However, since all our HRCs are continuous (with the exception of decelerations), we really want PDFs and not histograms. One approach for making PDFs is to assume an underlying distribution and then calculate point estimators for the parameters of that distribution using the data set. Unfortunately, we do not believe we can model our HRCs with any known distribution. Therefore, we resort to nonparametric methods for creating the PDFs. The method we use is called kernel density estimation. We will explain how this method works, the results of this method, and the conclusions from our PDFs in the following sections.

5.1 Kernel Density Estimation Definition

All our PDFs in this report are generated using kernel density estimation. In this section, we will explain how kernel density estimation works. Accordingly, first assume we have an independent and identically distributed sample x_1, x_2, \dots, x_n . The kernel density estimator $\hat{f}_X(x)$ for the PDF of X is then

$$\hat{f}_X(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x - x_i}{h}\right) \quad (5.1)$$

The function $K(\cdot)$ in Equation 5.1 is the kernel. For our report, we use the Epanechnikov kernel which is given below:

$$K(u) = \begin{cases} \frac{3}{4}(1 - u^2), & |u| \leq 1 \\ 0, & otherwise \end{cases} \quad (5.2)$$

One can think of the kernel density estimator (KDE) as a smoothed histogram. Consider Figure 5.1: the left panel shows a histogram for a data set, and the right panel shows the KDE for that same data set. In this data set, there are six data points represented by the ticks on the x-axis. The kernels in red are then calculate for each data point. Summing all the kernels together yields the solid blue KDE curve. Note that the author of this figure uses a normal kernel [10].

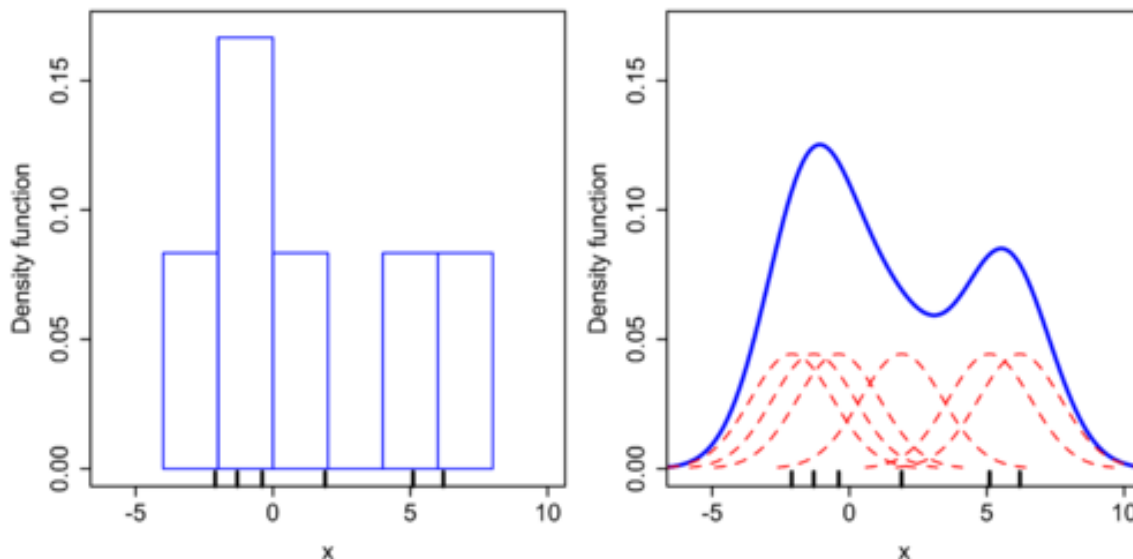


Figure 5.1: Comparison of Histogram and Corresponding KDE [10]

The kernel density estimation method is implemented in Matlab using the function `ksdensity()`. The two important choices we have to make are the kernel $K(\cdot)$ and the bandwidth h . We chose the Epanechnikov kernel because it is “optimal in the mean square error sense” [10]. The parameter h acts a smoothing parameter. Typically, we use Matlab’s default bandwidth. However, for reasons we will explain later, sometimes we must set the bandwidth ourselves. We use the Freedman-Diaconis method for choosing the bin width of a histogram to generate the bandwidth of our KDEs [3]. This method is given below:

$$bin\ width = 2\ IQR\ n^{-\frac{1}{3}} \quad (5.3)$$

Here, IQR stands for interquartile range and n is the number of observations. While Matlab uses a similar procedure to calculate the bandwidth, we may want to use a uniform bin width across various PDFs where the number of observations differ. Typically, we use the n and IQR from the unventilated organism 3 group because this is the smallest subgroup we analyze. Therefore, whenever we compare subgroups, we are limited by the observations of our smallest category.

5.2 Univariate PDF Figures

Univariate PDF Figures: `multiple_univariate_pdf_figures` →
`one_univariate_pdf_figure`

We can now generate univariate PDF figures using the kernel density estimation method. Figure 5.2 depicts an example of one of these figures. In Figure 5.2, we plot the raw HRCs for the group of sick half hours across all babies and the group of healthy half hours across all babies. Note that the healthy PDFs were generated using more half hours than the sick PDFs as indicated by the legend in the top left corner.

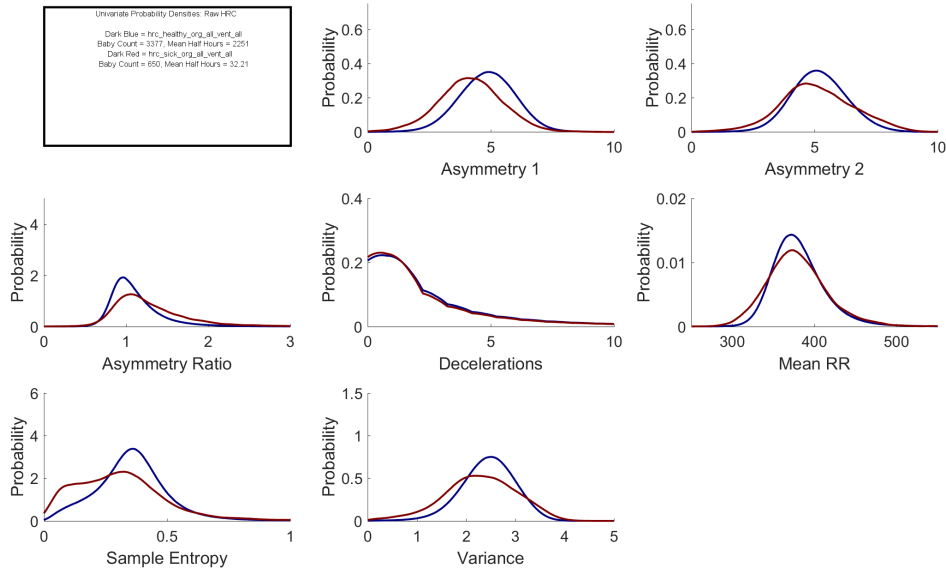


Figure 5.2: Univariate PDF Example

Appendix A contains more univariate PDF figures. Figure A.4 shows the PDFs for the raw HRCs from all the unventilated sick half hours separated by each organism. We also plot the PDF for the raw HRCs across all half hours for comparison. The next five figures correspond to the group of half hours from babies who experienced a septic event caused by organism 1 while they were in the NICU. There are five figures for this group for the five HRC subcategories (raw HRC, 10th percentile, 50th percentile, 90th percentile, and slope). Thus, we collectively have 35 subplots amongst the five figures for our 35 HRCs. On each subplot of each figure, we plot a PDF for the subgroups within the main group of organism 1 of sick ventilated, sick unventilated, healthy ventilated, and healthy unventilated. We can then repeat this process for every other organism. We can also repeat this process for any groups of half hours we wish to compare. However, the litany of additional PDF figures have been omitted from this report. For more information on the creation of these figures, see the documentation of `multiple_univariate_pdf_figures.m`.

5.3 Univariate PDF Discussion

When comparing healthy PDFs to sick PDFs using all relevant figures, we notice significant difference across all HRCs and HRC subcategories with the exception of decelerations. Consequently, we have our first major piece of support that the majority of the HRCs can distinguish between healthy and sick half hours. We also see significant differences when

comparing ventilated PDFs with unventilated PDFs for figures that display this information. This result leads us to separate the analysis of ventilated and unventilated half hours in the later chapters. Another result we find is that there is little difference in the PDFs of the five organisms. We thus conclude that it is much more difficult to determine the type of organism that caused the illness if we know a baby is ill when compared to determining if a baby is sick or healthy. At this point in the report, we will stop separating half hours by organism. Rather, we will focus on separating groups by healthy or sick and ventilated or unventilated.

Chapter 6

Univariate Risks

In the previous section, we plotted the PDFs of two groups in order to compare them. While this method provides a good visual way of comparing two groups, we cannot take away any quantitative information. Ideally, we want to define a measure, which we will call the risk, that relates the two PDFs. In this section, we will define and analyze the risk.

6.1 Risk Definition

The risk gives us a number indicating if a half hour is more likely to be in the sick group or healthy group. Accordingly, let us first define the probability of a half hour being sick or $P(sick)$:

$$P(sick) = \frac{\text{number of sick half hours}}{\text{number of sick half hours} + \text{number of healthy half hours}} \quad (6.1)$$

Note that the number of sick/healthy half hours changes depending on the subgroup we are concerned with. For example, if we are calculating a risk for unventilated half hours, the number of sick/healthy half hours must be calculated using only this group. Next, we will define the probability of a half hour being sick given an HRC signal or $P(sick|signal)$:

$$P(sick|signal) = \frac{P(signal|sick)P(sick)}{P(signal|healthy)P(healthy) + P(signal|sick)P(sick)} \quad (6.2)$$

Since we are dealing with univariate risks, the variable “signal” will refer to one HRC (e.g. $P(sick|variance = x)$). When we discuss bivariate risks in Chapter 8, signal will refer to HRC 1 and HRC 2 (e.g. $P(sick|variance = x \text{ and sample entropy} = y)$). Next, we can calculate the terms $P(signal|sick)$ and $P(signal|healthy)$ using the PDFs generated in Chapter 5. Because those terms are probabilities, we integrate the PDFs over the signal’s respective bin width centered around the value of the signal to get those two terms. Moreover, integrating the PDF allows us to find the probability between two points. We therefore estimate to probability terms by integrating over a small area defined by the bin width. Finally, our definition of risk is

$$Risk = \ln \left(\frac{P(sick|signal)}{P(sick)} \right) \quad (6.3)$$

Since we take the natural logarithm of this fraction, any number above 0 indicates that the half hour is more likely to be sick given the signal. In the next section, we will create figures that plot the risk over typical values of each HRC signal.

6.2 Univariate Risk Figures

Univariate Risk Figures: `multiple_univariate_risk_figures` →
`one_univariate_risk_figure` → `one_risk_matrix` → `one_prob_matrix`

Figure 6.1 gives an example of what one of the univariate risk figures looks like. This figure illustrates the risk of a half hour being sick across all half hours (ventilated and unventilated). The figure shows seven subplots for the seven HRCs. On each subplot, we have five curves for the five HRC subcategories. Accordingly, this figure gives the risk for all 35 HRCs. If the probability of a certain HRC is very small, we do not plot the risk. For example, a variance value of 4.5 is so rare that we cut off the curve before this point.

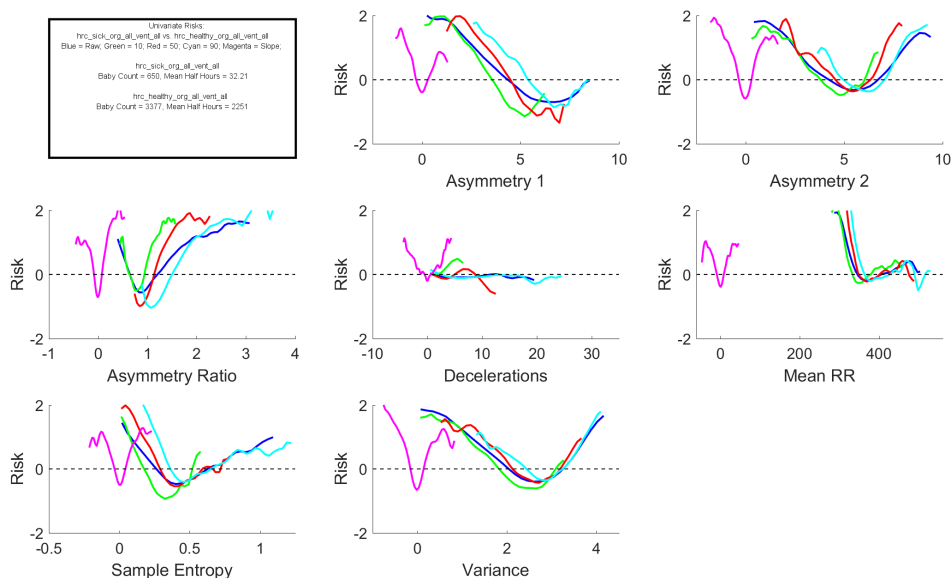


Figure 6.1: Univariate Risk Example

Appendix A contains two more risk figures: one for only unventilated half hours and another for only ventilated half hours. Again, we are not separating the half hours by organism anymore. However, if we were to separate half hours this way, the risk of organism x compared to organism y would be very flat and around 0 across all HRCs, indicating that the HRCs do not do a very good job at distinguishing between two organisms. For more information about how we create the univariate risk figures, please see the documentation of `multiple_univariate_risk_figures.m`.

6.3 Univariate Risk Discussion

When looking at the univariate risk figures, we hope to find curves that are significantly

above the dotted 0 line in each subplot. Effectively, such a scenario would indicate that these HRC values would be very predictive of sepsis. However, the figures can be misleading. For example, consider a variance value of 4 in Figure 6.1. From the figure, a half hour with a variance of 4 would indicate that the half hour is e^2 times more likely than the typical half hour to be sick. We would thus conclude that variance is a very strong predictor of sepsis. The problem with this analysis is that a variance of 4 is very unlikely, so such a value holds little predictive power. What one really must do is look at the univariate risk figures with the univariate PDFs of each HRC in mind. The most useful HRCs would then be ones whose risk is relatively large in areas where the sick PDF is also large. To reiterate this point one more time, if the $P(\text{sick}|\text{variance} = 4) = 0.002$ and $P(\text{healthy}|\text{variance} = 4) = 0.001$, we would conclude a baby is e^2 times as more likely to be sick than healthy when the variance is 4. However, since the sick probability is so small, the risk is not extremely useful.

If we ignore the distributions of the HRCs for a moment, Figure 6.1 demonstrates that the following values are predictive of sepsis: low asymmetry 1 values, low and high asymmetry 2 values, high asymmetry ratio values, low mean RR values, low sample entropy values, and low and high variance values. Furthermore, low and high values for all slopes seem to predict sepsis as well.

Ideally, we would plot the sick and healthy PDFs of each HRC on the risk plots. However, with five risks on each subplot, adding ten more PDFs would become quite cumbersome. Fortunately, the figures in Chapter 9 will give a similar illustration of such a figure.

Chapter 7

Bivariate Probability Density Functions

After developing methods to calculate the probability of illness given one HRC, we can extend these methods to two HRCs. We will still use the kernel density estimation method to generate PDFs. However, now we will sum over bivariate Epanechnikov kernels instead of univariate ones. Because Matlab's default bandwidth smooths the bivariate PDFs too much, we set our own bandwidths for the bivariate PDFs. We use bandwidths that are half of the bin widths calculated by Equation 5.3. While our choice for halving the bandwidths is empirical, the bandwidths of the HRCs for a bivariate PDF must be smaller than the bandwidths for a univariate PDF so that the rectangle we smooth over in the bivariate case is somewhat proportional to the line segment we smooth over in the univariate case.

Again, we use Matlab's `ksdensity()` function to generate the bivariate KDEs. The rest of the chapter will describe the resulting bivariate figures and analyze their meaning.

7.1 Bivariate PDF Figures

Bivariate PDF Figures: `multiple_bivariate_pdf_figures` →
`one_bivariate_pdf_figure`

Similar to the univariate case, we will create bivariate PDFs for different combinations of HRCs. However, with 35 HRCs, we would end up with 595 bivariate PDFs. Therefore, we will only analyze the six bivariate PDFs which we empirically flagged as the most useful: raw variance vs. raw sample entropy, raw variance vs. raw asymmetry ratio, raw variance vs. raw decelerations, raw sample entropy vs. raw asymmetry ratio, raw sample entropy vs. raw decelerations, and raw asymmetry ratio vs. raw deceleration.

Figure 7.1 provides an example of a figure with these six bivariate PDFs. The figure compares the sick half hours across all babies against the healthy half hours across all babies. Appendix A provides two more examples of these figures: one comparing unventilated sick and healthy half hours and one comparing ventilated sick and healthy half hours. The documentation of `multiple_bivariate_pdf_figures.m` contains more info on these figures.

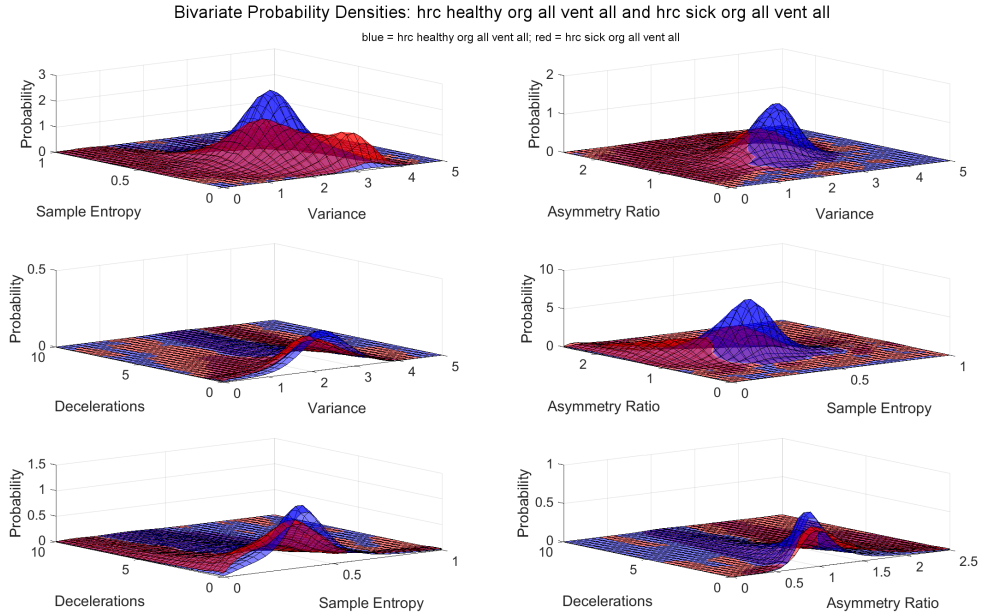


Figure 7.1: Bivariate PDF Example

7.2 Bivariate PDF Discussion

Parallel to the analysis for the univariate PDFs, we are looking for disparate sick and healthy bivariate PDFs for one set of HRCs. Figure 7.1 illustrates that the combination of variance and sample entropy is a very good predictor of sepsis. Any set with deceleration, alternatively, produces very similar sick and healthy PDFs. Overall, Figure 7.1 shows strong support for the combinations of raw variance with raw sample entropy, raw variance with raw asymmetry ratio, and raw sample entropy with raw asymmetry as good predictors. Other HRCs, however, might be more useful for different groups and needs to be looked into further.

Theoretically, with an infinite number of observations, we could calculate the probability of illness given n HRCs using an n -variate PDF. However, we believe with our data set, the best we will be able to calculate is probabilities using two HRCs. Visualizing the probabilities given more than two HRCs would also be quite a challenge. Therefore, we will have to resort to different statistical techniques in the later chapters which will allow us to incorporate more than two HRCs. However, first we will analyze the bivariate risks in the next chapter.

Chapter 8

Bivariate Risks

The bivariate risks mirror Equation 6.3 except this time *signal* refers to two HRC signals. Furthermore, we calculate the terms $P(\text{signal}|\text{healthy})$ and $P(\text{signal}|\text{sick})$ by integrating the bivariate PDFs over a small rectangle around the value in question. The dimensions of the rectangle are given by the bin widths of the two HRCs. The next step is to make bivariate risk figures similar to the univariate ones.

8.1 Bivariate Risk Figures

Bivariate Risk Figures: `multiple_bivariate_risk_figures` → `one_bivariate_risk_figure` → `one_bivariate_risk_plot` → `one_risk_matrix` → `one_prob_matrix`

Figure 8.1 depicts the bivariate risks corresponding to the bivariate PDFs in Figure 7.1. Furthermore, Appendix A gives the bivariate risk figures for the bivariate PDF figures in the appendix. The function `multiple_bivariate_risk_figures.m` creates these figures.

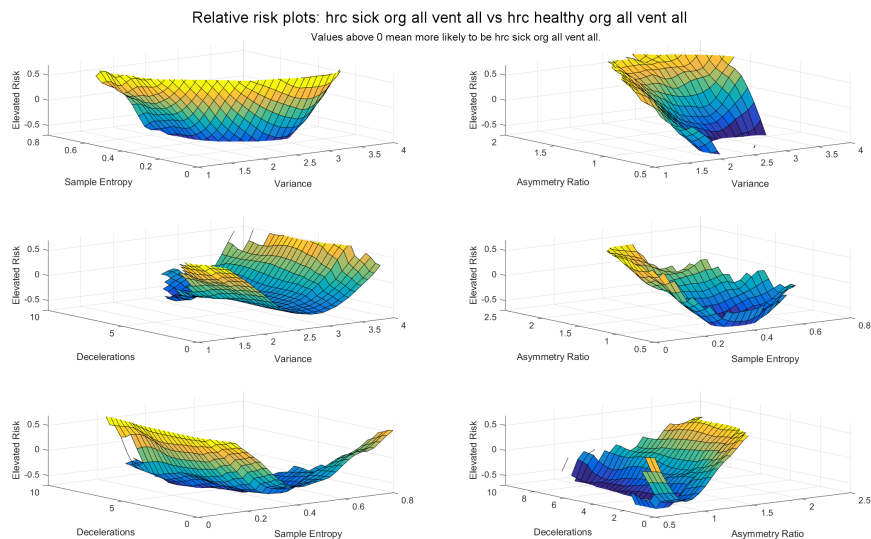


Figure 8.1: Bivariate Risk Example

8.2 Bivariate Risk Discussion

As with the univariate risks, one must keep the corresponding bivariate PDFs in mind when analyzing the bivariate risks. While high risks show the HRC set is useful at predicting sepsis, the set's predictive power is diminished if there is limited data in the high risk area. Any risk above the $z = 0$ plane in Figure 8.1 indicates the half hour is more likely to be sick given that set of HRCs. The analysis of Figure 8.1 provides the same conclusions that we found in Chapter 7, *i.e.*, low and high variance are useful, low sample entropy is useful, etc.

Chapter 9

Single Variable Logistic Regression

As stated at the end of Chapter 7, we cannot simply continue calculating n -variate PDFs. Instead, we will resort to logistic regression in order to calculate the probability of illness given n HRCs. Before we can build up to logistic regression with n HRCs, we will start with a single HRC. The rest of this section details the mathematics of single variable logistic regression as well as the resulting figures and analysis.

9.1 Single Variable Logistic Regression Definition

Consider a response vector $[y_1, y_2, \dots, y_n]^T$ drawn from a binomial distribution with support $y_i = \{0, 1\}$. Now consider a set of explanatory variables $[x_1, x_2, \dots, x_q]^T$. Using a logistic regression model, the probability that $Y_i = 1$ is then

$$P(Y_i = 1) = \frac{e^{\beta \mathbf{x}}}{e^{\beta \mathbf{x}} + 1} \quad (9.1)$$

where \mathbf{x} is the vector of explanatory variables and β is the vector of coefficients for the explanatory variables. The coefficients are calculated using a least squares method which requires the response vector and the matrix of associate explanatory variables. Moreover, each response variable y_i has an associated \mathbf{x}_i for a total of n \mathbf{x} -vectors. The i^{th} \mathbf{x} -vector becomes one row in the matrix of explanatory variables [11].

For our single variable logistic regression models, \mathbf{x} will have two or three terms depending on whether we use a linear or quadratic model. In order to explain the terms in the vector, consider the example of predicting sepsis using the single HRC of raw variance. The two possibilities for \mathbf{x} are given below:

$$\text{Linear} : \mathbf{x} = [1 \text{ variance}]^T$$

$$\text{Quadratic} : \mathbf{x} = [1 \text{ variance variance}^2]^T$$

Note that the 1 in each vector represents the y -intercept. The term *variance* is simply the raw variance for that half hour and the term *variance*² is the square of the raw variance.

For this report, our response vector is 1 if the half hour is sick and 0 if the half hour is healthy. We can then use the Matlab function `fitglm()` to calculate the coefficients for our model using the response vector and the matrix of explanatory variables, i.e., the HRCs. Next, we will construct figures for the probability of illness given various values of a single HRC.

9.2 Single Variable Logistic Regression Figures

Single Variable Logistic Regression Figures: `multiple_univariate_logistic_figures` → `one_univariate_logistic_figure` → `one_risk_matrix` → `one_prob_matrix`

We can now run our single variable logistic regression on all 35 HRCs. Figure 9.1 gives the results of the regression for the raw HRCs using all the healthy and all the sick half hours. The logistic regression probabilities are given by the blue curve. We use a quadratic model for each of these figures. The red curve on each subplot, which we call the Bayesian probability, is simply the $P(\text{sick}|\text{signal})$ term from the univariate risks. Finally, the red and blue areas are the distributions of the HRCs where the blue region represents healthy half hours and the red region represents the sick half hours. Note that the shaded regions are stacked on top of one another, i.e., no region is being hidden. One can interpret the entire shaded regions as the joint distribution of healthy and sick half hours. The scale for the joint distribution is given by the right y-axis.

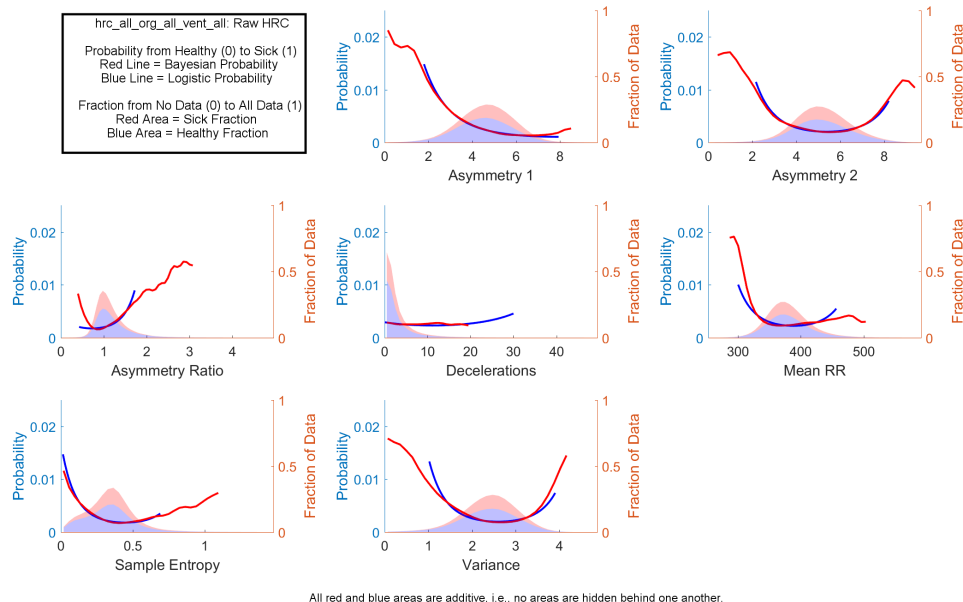


Figure 9.1: Single Variable Logistic Regression Example

Appendix A gives the corresponding figure for the 10th percentile, 50th percentile, 90th percentile, and Slope HRCs. Furthermore, one could calculate the logistic probabilities for all 35 HRCs using a different group of healthy and sick half hours such as only the ventilated or only the unventilated half hours. For more information on these figures, please see the documentation of `multiple_univariate_logistic_figures.m`.

9.3 Single Variable Logistic Regression Discussion

The analysis of the single variable logistic regression follows closely to the analysis the univariate risks. First, recall that the red curves are the $P(\text{sick}|\text{signal})$ terms from the univariate risk calculations. One can then think of the red curves as the empirical results

and the blue curves as our fitted model. Likewise, as we addressed for the univariate risks, the most predictive HRCs are the ones who have large sick probabilities/risks in areas with large sick PDFs. The PDFs on the logistic figures, therefore, serve a very useful purpose.

We start with single variable logistic regression to first check our empirical results. We also start with one variable to try to cut down on the number of HRCs before we move to the n -variable logistic regression. If we use all 35 HRCs for the logistic regression, we would have 666 terms in our \mathbf{x} -vector for a quadratic model. However, a model with so many terms would severely over-fit our data. The single variable logistic regression can, therefore, provide one method to hand pick a few HRCs we think are the most predictive.

Looking at Figure 9.1, we see a lot of HRCs with high probabilities for low values and high values of the HRCs. This phenomenon gives us strong evidence that we should use a quadratic model over a linear model. For example, the variance subplot in Figure 9.1 shows high probabilities on both tails. A linear fit would only capture one extreme. Therefore, we need a reaction term, which gives us a quadratic model, to capture both sides of the distribution.

Not shown in this report are the single variable logistic figure for only the ventilated group and only the unventilated group. The important result from these sets of figures is that the ventilated probabilities whether sick or healthy are about six times greater than the unventilated probabilities. From this finding, one might want to include ventilation status as an additional HRC. However, we do not include ventilation status because different NICUs might have different standards for ventilating babies. Despite this concern, babies who are ventilated do seem to be more likely to be sick. Therefore, we continue to separate half hours into ventilated and unventilated because there may be different optimal parameters and thresholds for each group when determining if a baby is sick or healthy.

Chapter 10

HeRO Score

The groundwork has now been set to conduct n -variable logistic regression which we call the HeRO score. This section will describe various HeRO score models we have developed and the findings from these models.

10.1 HeRO Score Definition

In general, the HeRO score is an n -variable logistic regression model. The formula for the HeRO score is given below:

$$\text{HeRO Score} = \frac{e^{\beta\mathbf{x}}}{e^{\beta\mathbf{x}} + 1} \cdot \frac{1}{\mu_0} \quad (10.1)$$

In this formula, we multiply the logistic regression model by the parameter μ_0^{-1} where μ_0 represents the probability of a half hour being sick. Multiplying by this coefficient means that a HeRO score of 1 indicates that the half hour is just as likely to be sick as any half hour. A HeRO score of 2 would mean the half hour is twice as likely to be sick compared to the typical half hour.

Our explanatory variables vector \mathbf{x} can take on many forms. Below gives three models we will analyze in this section:

$$\text{Hrch} : \mathbf{x} = [1 \text{ var}_{10} (\text{asym2} - \text{asym1})_{10}]^T$$

$$\text{Hrcg} : \mathbf{x} = [1 \text{ var}_{50} \text{ sampen}_{10} \text{ asym1}_{50} \text{ asym2}_{50}]^T$$

$$\begin{aligned} \text{Dienstman} : \mathbf{x} = & [1 \text{ asym1}_{10} \text{ asym1}_{\text{slope}} \text{ asym2}_{90} \text{ asym2}_{\text{slope}} \text{ asym_ratio}_{50} \text{ asym_ratio}_{\text{slope}} \\ & \text{decls}_{90} \text{ sampen}_{10} \text{ sampen}_{\text{slope}} \text{ var}_{10} \text{ var}_{90} \text{ var}_{\text{slope}} \\ & \text{asym1}_{10} * \text{asym1}_{\text{slope}} \text{ asym1}_{10} * \text{asym2}_{90} \dots \text{asym1}_{10} * \text{var}_{\text{slope}} \dots \\ & \text{asym1}_{10}^2 \dots \text{var}_{\text{slope}}^2]^T \end{aligned}$$

Note that the *Hrch* and *Hrcg* models are linear models while the *Dienstman* model is a quadratic model. The *Hrch* and *Hrcg* model were developed by Prof. Douglas Lake while the *Dienstman* model was developed for this report. As a result, the coefficients for the *Hrch* and *Hrcg* model were calculated using a different data set and with a slightly different implementation than the coefficients for the *Dienstman* model. In total, the *Dienstman* model contains 91 coefficients and contains every reaction term of the 12 initial HRCs. We

picked the HRCs for the *Dienstman* model based on the single variable logistic probabilities that we thought to be the most predictive. We tried not to include too many HRCs for fear of over-fitting. For more information about how to calculate the coefficients for the *Dienstman* model, please see the documentation of `csv_logistic_coeffs.m`.

All HeRO scores use percentiles and slopes instead of the raw HRCs in order to reduce noise when we plot a time series of the HeRO scores in the next section. If HeRO scores only used information from the current half hour, then one outlier would cause the HeRO score to fluctuate considerably, and because the babies are frequently disturbed by the nurses, we can get a significant number of outliers. Therefore, if we look at var_{10} for example, the second lowest variance value over the previous 24 half hours would go into the HeRO score and not the lowest, reducing outliers and noise (10th percentile usually translates to the second lowest value for a 24 half hour window).

10.2 HeRO Score Figures

HeRO Score Figures: `multiple_hero_score_figures` → `one_hero_score_figure`

After we calculate the coefficients for all the HeRO score models, we can make time series of the HeRO scores for every baby who had an event as we did for the time series in Chapter 4. Figure 10.2 provides an example of one of these figures. The vertical black line represents the time of the septic event and the four curves represent four different HeRO scores. The green line, which we call the legacy HeRO score or just HeRO score, is the maximum between the Hrch and Hrcg models, which are indicated on the figure by the plus makers and circle markers, respectively. The gray line is the *Dienstman* HeRO score discussed previously. Finally, the *Dienstman Vent* and *Dienstman Nonvent* HeRO scores are the *Dienstman* model, but instead of using all the half hours to calculate the coefficients, we only use the ventilated or unventilated half hours. We plot *Dienstman Vent* on the figure only when the baby is ventilated, and otherwise we plot *Dienstman Nonvent*. Appendix A provides more examples of these figures for other events. All HeRO score figures were developed using `multiple_hero_score_figures.m`.

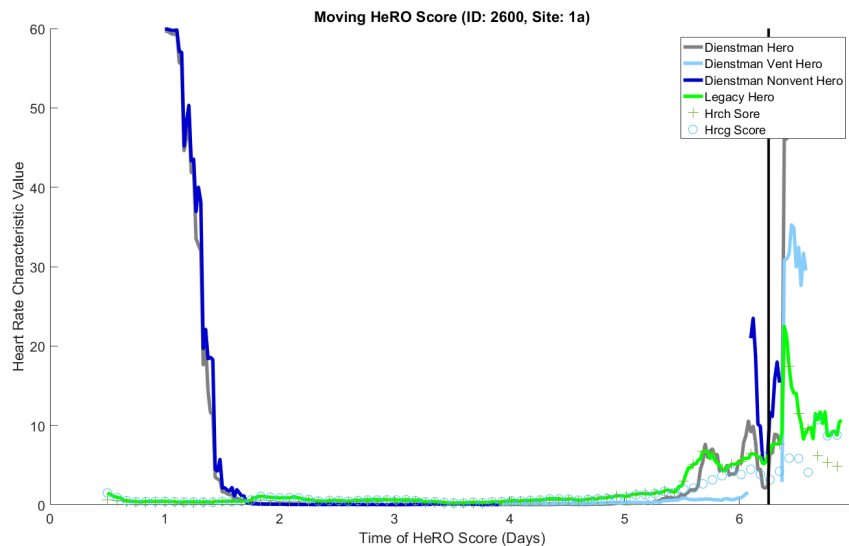


Figure 10.1: HeRO Score Figure Example

After we have created these figures for all events, we can take the average HeRO score leading up to the event in the same manner as we took the average HRC for the time series figures. Figure 10.2 presents the resulting figure for the average HeRO score. Recall that not every baby will have data for every half hour leading up to an event. Subsequently, we record how many half hours we use for each half hour relative to the event and record the average of these numbers in the title of the figure. It should also be noted that *Dienstman Vent* and *Dienstman Nonvent* use about half as many observations as the number noted in the title since a baby can never be ventilated and unventilated at the same time.

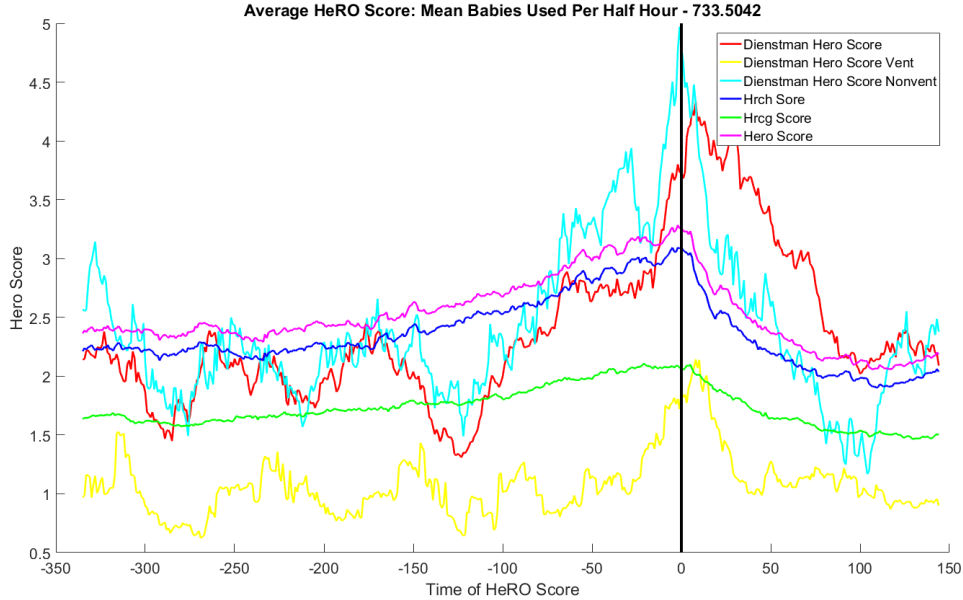


Figure 10.2: Average HeRO Score Figure

10.3 HeRO Score Discussion

The first major feature of the individual HeRO figures is that the *Dienstman* model seems to over fit the data. HeRO scores are either extremely high or extremely low, which is a classic sign of over fitting. The individual figures are also very noisy and do not always provide a clear indication that a high HeRO score will indicate the baby is sick.

The average figure, on the other hand, does show us that the HeRO scores are predicting sepsis. Unlike the time series figures, the HeRO scores increase right before an event on average. Thus, while no individual HRC showed significant predictive power in the average time series, a combination of multiple HRCs certainly provides warning for the onset of sepsis. The legacy HeRO score also provides a very smooth average. This feature is extremely useful because it gives medical professionals more confidence that a high HeRO score is a true indication of illness and not a false positive. However, the individual legacy HeRO scores are still noisy, so the average in this case may be deceiving. We hypothesize that the noise for the average *Dienstman* models comes from over-fitting as one large score at half hour x will skew the average. However, one interesting feature is that the *Dienstman Vent* and *Dienstman Nonvent* scores are significantly different, giving us more support that we should treat the ventilated and unventilated half hours separately.

In general, we hope to see an increase in the HeRO leading up to the event. However, at present, we have not implemented any method for ranking the HeRO scores apart from just looking to see if they increase smoothly before the event. A gradual increase might be useful as it could give early warning that the baby is ill. Conversely, a sharp spike right before the event also has merit as it provides a clear threshold for the moment the baby requires attention. Accordingly, more work needs to be done in order to determine which model performs the best. This will be especially important for testing the performance of any future HeRO score against current ones.

Lastly, note that we used all the half hours to calculate the coefficients for the *Dienstman* scores. Ideally, we should separate the half hours into a learning set and testing set. We would then calculate the coefficients using only the learning set and produce figures for the testing set based off these coefficients. Since we did not take this step, this area could be another cause for noise and over-fitting. Also note that we have many more healthy half hours than sick half hours. Thus, it might also be necessary to randomly select a subgroup of the healthy half hours so we use the same number of healthy and sick half hours to calculate the coefficients. Accordingly, we need to look into this issue more carefully.

Chapter 11

Conclusion

After rigorously analyzing the predictive power of our 35 HRCs, there are many important results to summarize. Most of these findings also lead to more questions and further areas to investigate. We will address both topics in the following sections.

11.1 Discussion

Recall that one of our early goals was to try to distinguish between the types of invading organisms that cause illness. Our original hypothesis was that low variability, or variance, is predictive of gram-positive bacteria and decelerations are predictive of gram-negative bacteria. By analyzing the univariate PDFs, we saw that no HRC performed particularly well in identifying the invading organism. Therefore, we would either need to develop new HRCs to accomplish this goal or analyze the current HRCs using a different method.

We were, however, able to show substantial difference in the PDFs of sick and healthy half hours. The HRCs of variance, sample entropy, and asymmetry ratio were particularly useful in this respect. These three HRCs had univariate and bivariate PDFs that differed significantly between their respective sick and healthy PDFs. The univariate and bivariate risks also demonstrated this discovery.

We also repeatedly saw that ventilated and unventilated half hours behave very differently. The HRCs for these two groups gave very distinct distributions, risks, and HeRO scores. Subsequently, we have given strong evidence for the use of separate parameters, thresholds, and baselines for these two groups in any future work.

Lastly, the HeRO scores showed promising signs for the ability to predict illness using many HRCs. We are limited in the analysis we could do through Bayesian methods by the size of our data. Thus, a logistic regression based model proved to be a viable way to predict illness using multiple HRCs. However, the HeRO scores were very noisy for the individual HeRO score figures, and the *Dienstman* scores seemed to over-fit that data. Thus, we need to address these issues in future work.

11.2 Future Work

A significant portion of the future work involves improving the HeRO score. As mentioned at the end of Chapter 10, we need to reduce the noise in some of the scores. One possible solution is to reduce the number terms in the explanatory vector. We could remove some reaction terms so we have a mix of a pure linear and pure quadratic model. We could

also remove some HRCs entirely. We also need to separate data into a learning set and testing set as well as investigate whether we need to use the same number of healthy half hours as sick half hours for calculating the coefficients.

We also want to try to create the best HeRO score possible for predicting sepsis. Such a goal might require adding more HRCs to a model. We could also have HeRO scores for other groups beyond just ventilated and unventilated. These groups may be based on gestational age, post-menstrual age, gender, and race. The HeRO scores for these groups might also require different HRCs as opposed to just different coefficients for the same HRCs. We would also need a reasonable way to compare the performance of all these different HeRO scores in order to determine which one we would prefer the most. Such a test should consider how many false positives and negatives the score gives and how early of a warning does the score give.

Finally, we could also try to look towards other statistical techniques for predicting sepsis beyond logistic regression. Possibly models include support vector machines or neural networks.

Appendix A

Additional Figures

This appendix contains additional figures not in the chapters above. Note that this appendix is not a complete list of figures produced from the research. Because of the large number of figures, we have elected to only give a few more examples of the types of figures we discussed. However, we can produce more figures upon request.

A.1 Time Series Figures

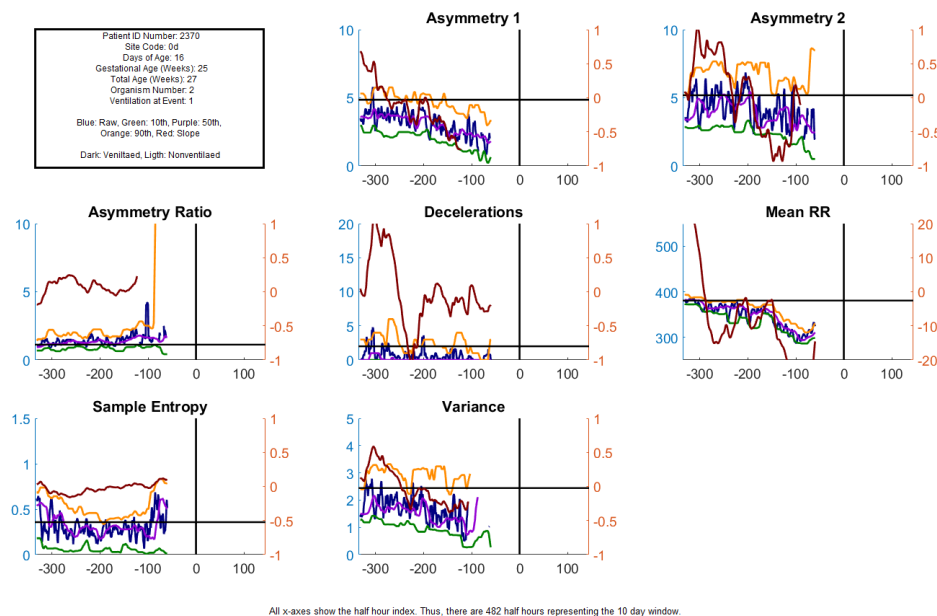


Figure A.1

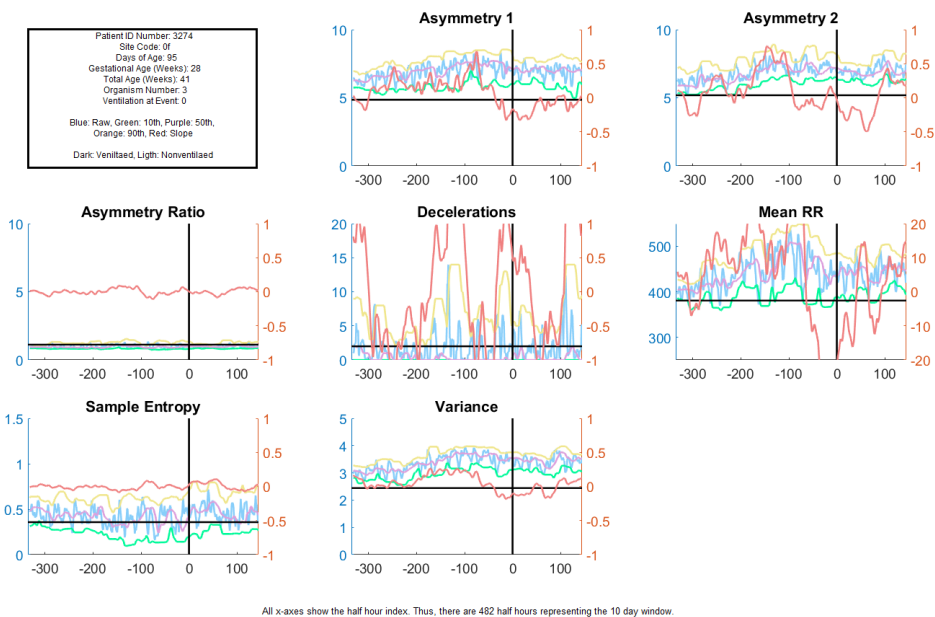


Figure A.2

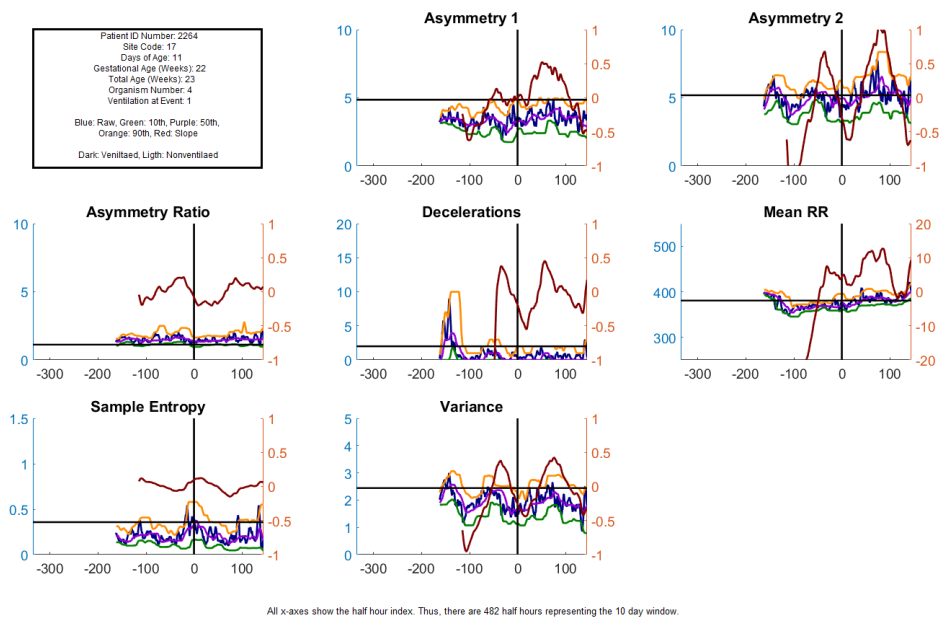


Figure A.3

A.2 Univariate PDF Figures

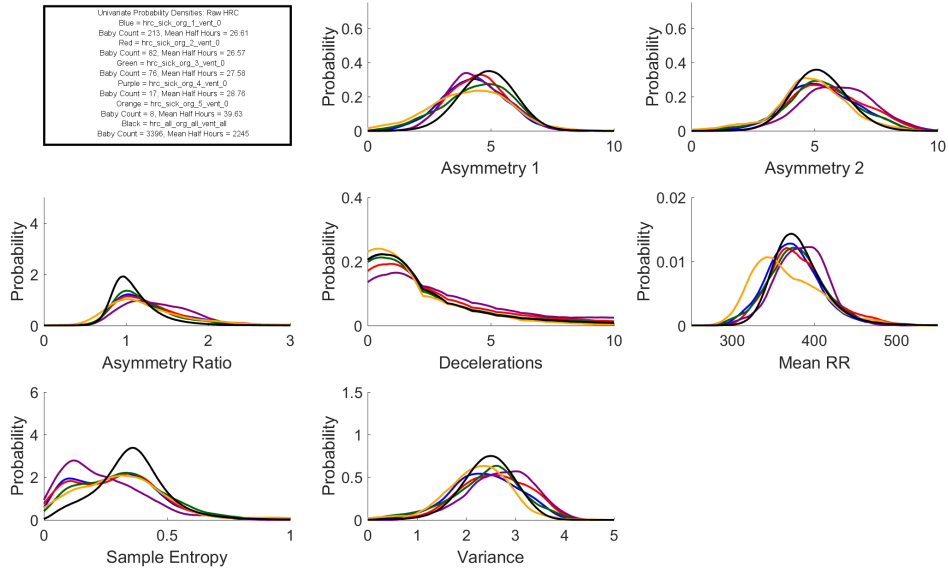


Figure A.4

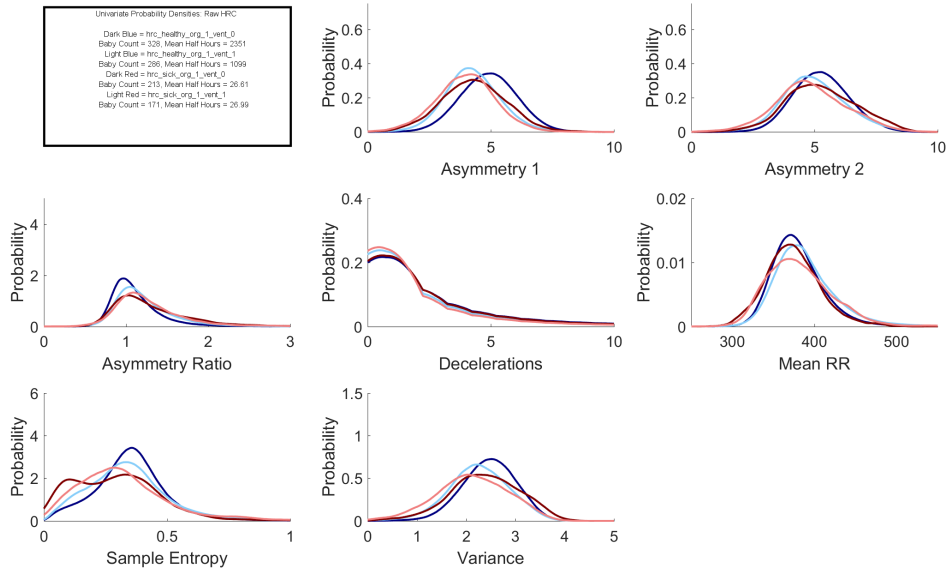


Figure A.5

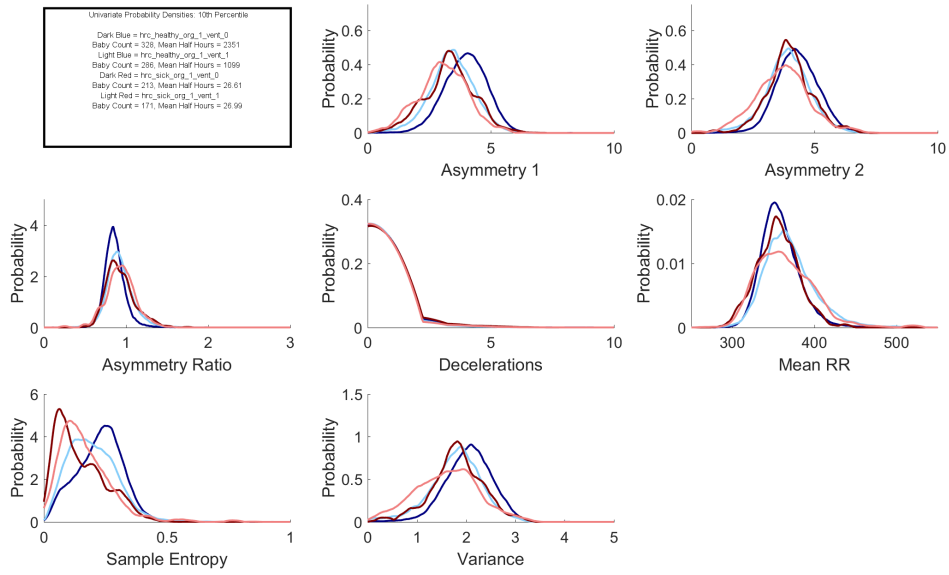


Figure A.6

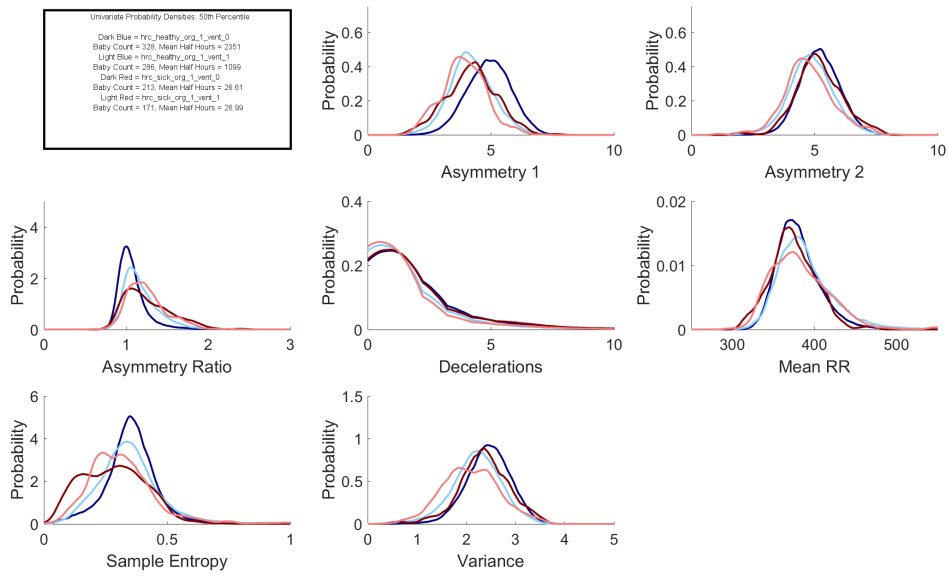


Figure A.7

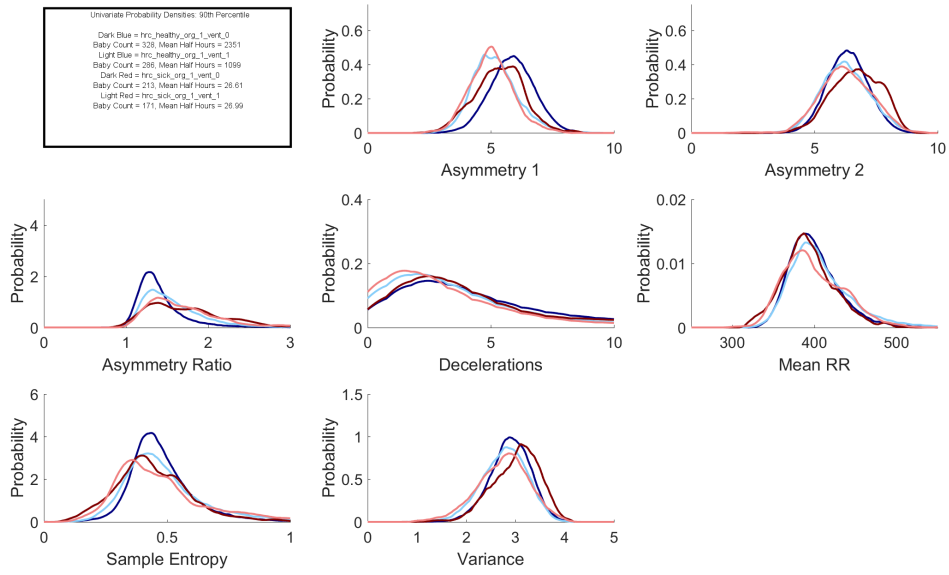


Figure A.8: Univariate PDF 5

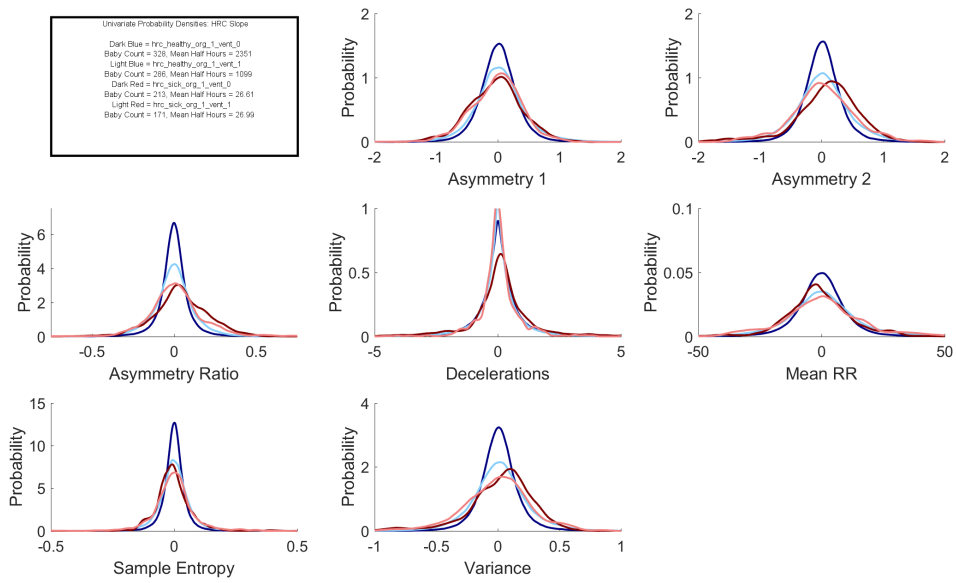


Figure A.9

A.3 Univariate Risk Figures

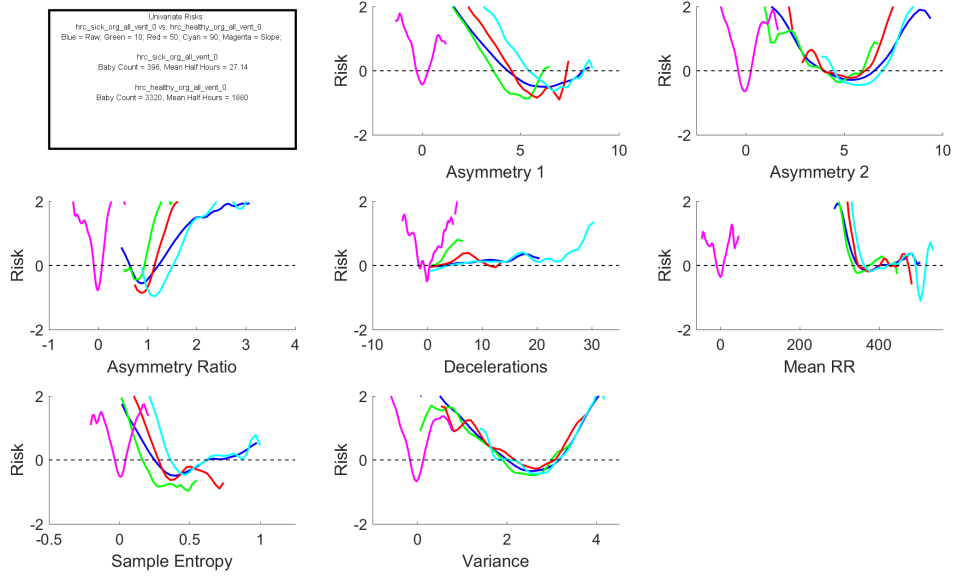


Figure A.10

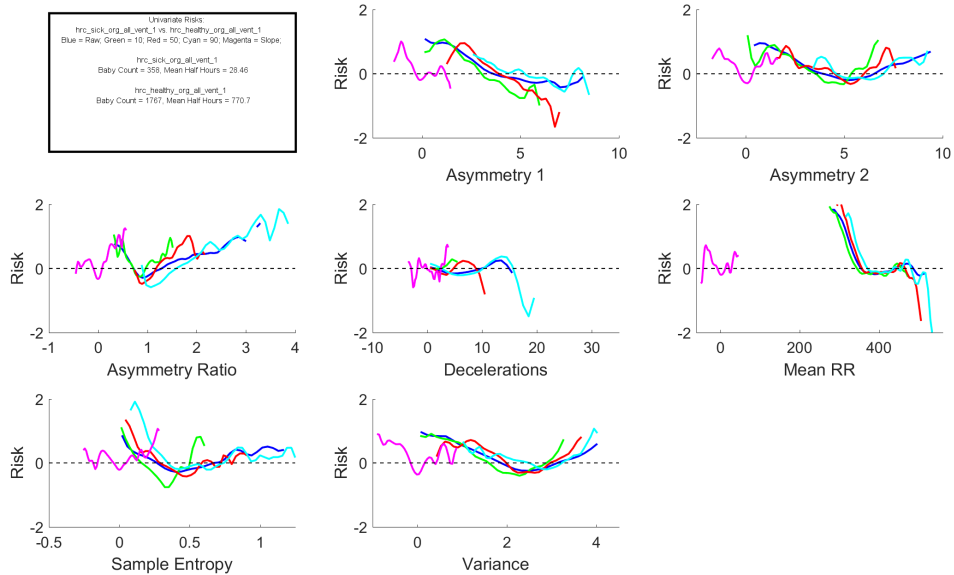


Figure A.11

A.4 Bivariate PDF Figures

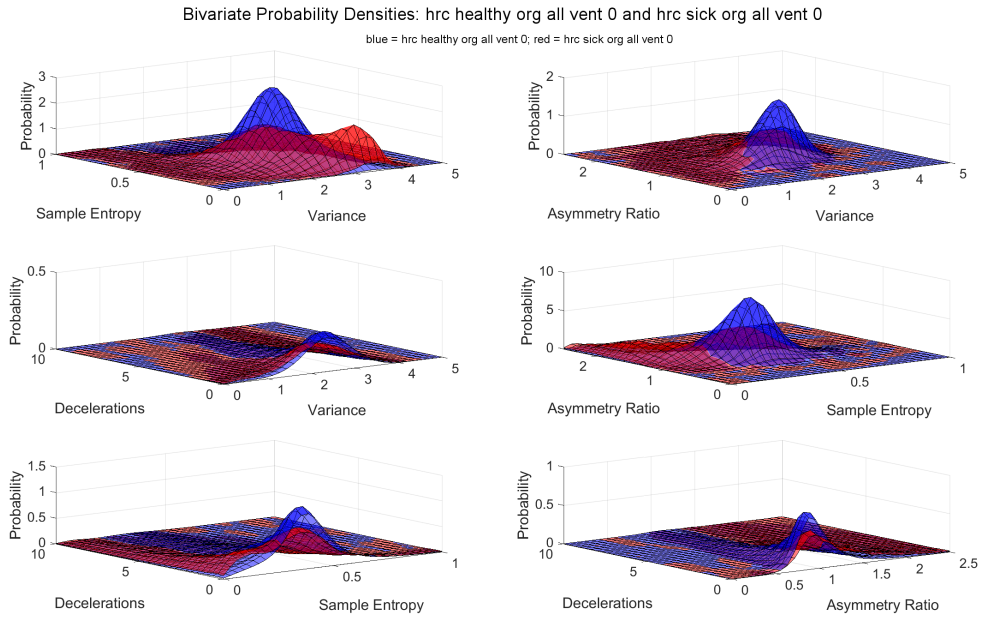


Figure A.12

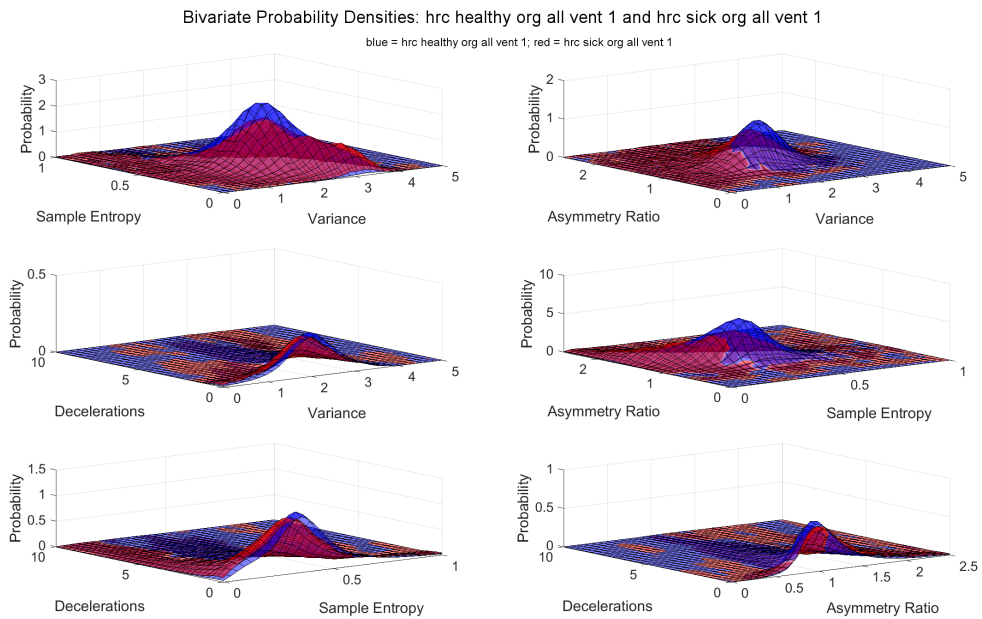


Figure A.13

A.5 Bivariate Risk Figures

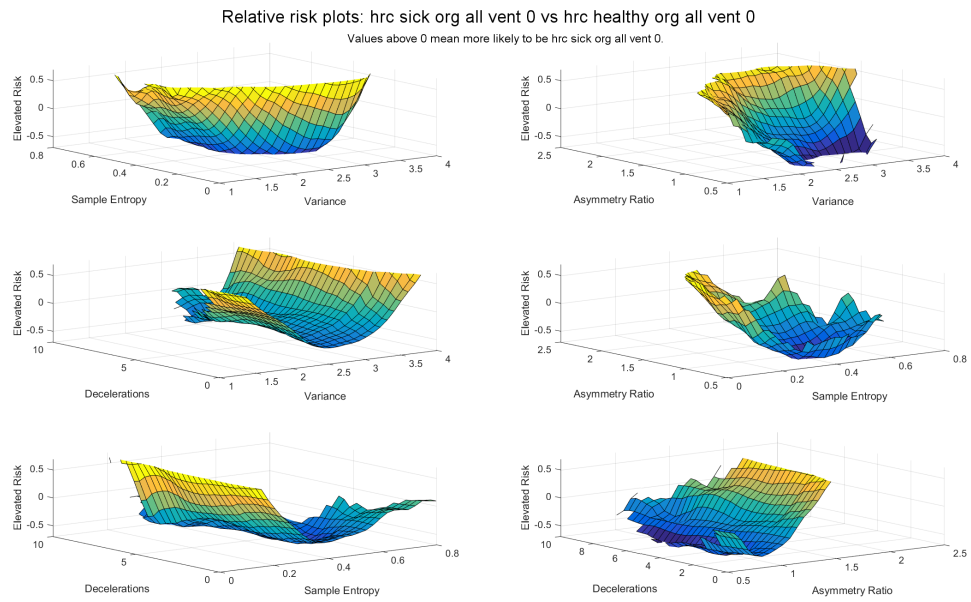


Figure A.14

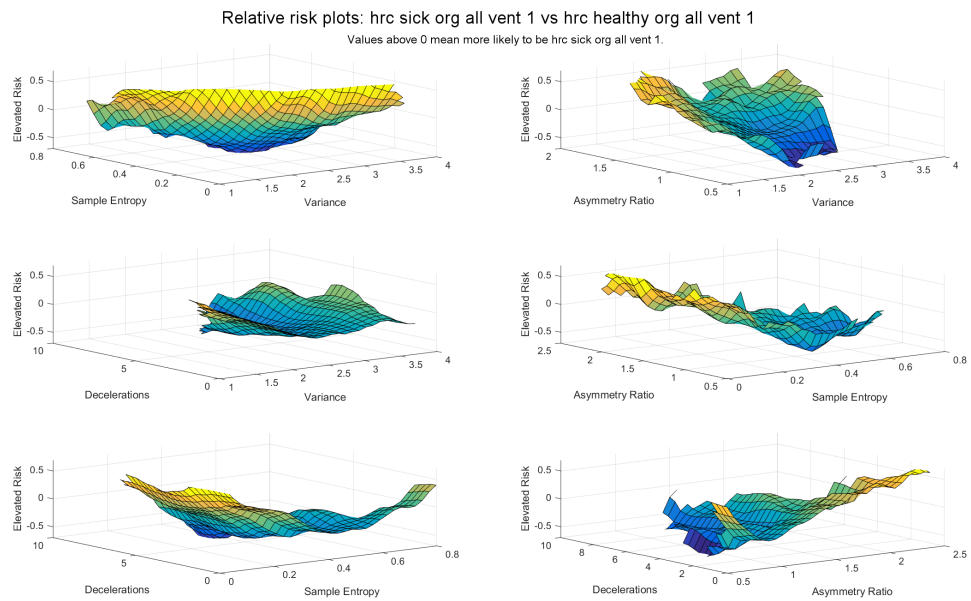


Figure A.15

A.6 Single Variable Logistic Figures

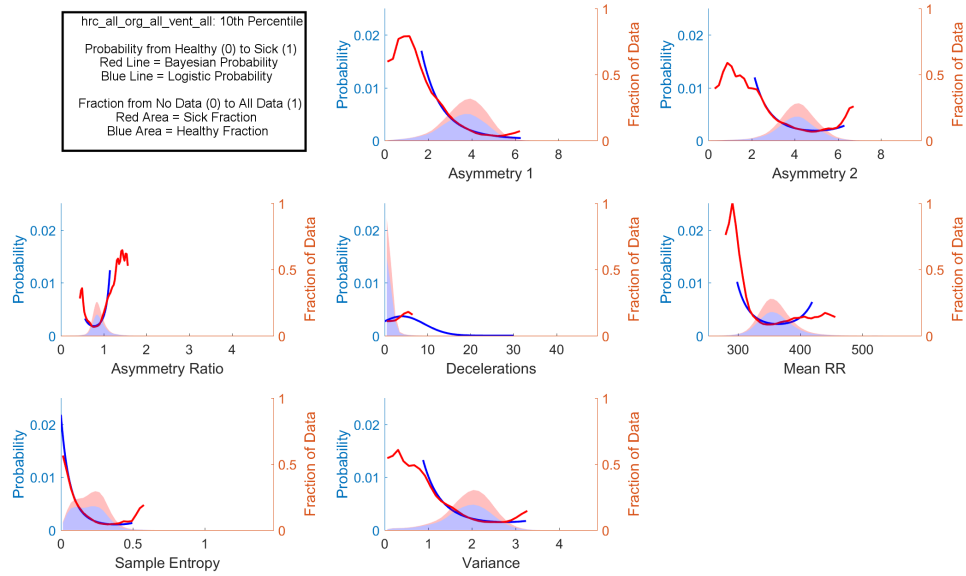


Figure A.16

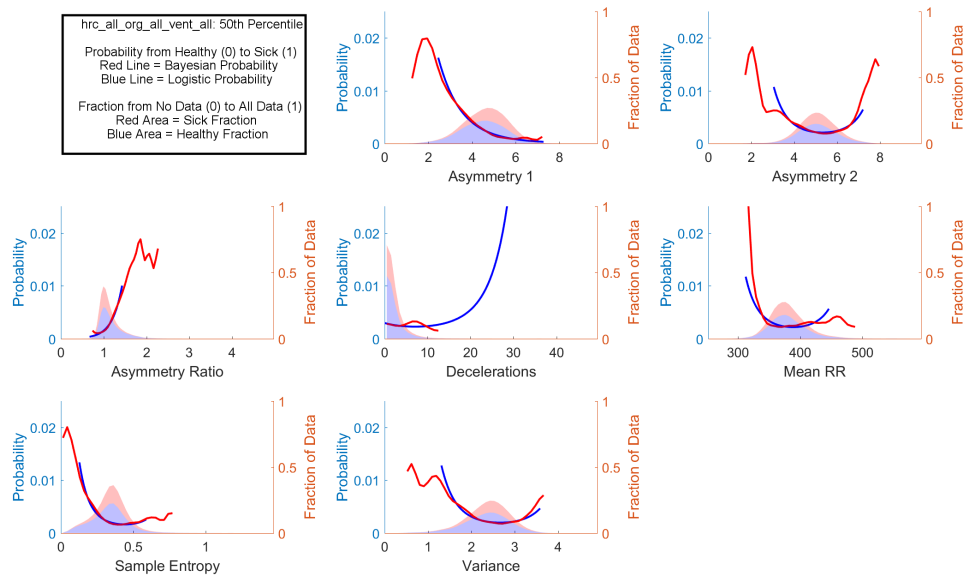
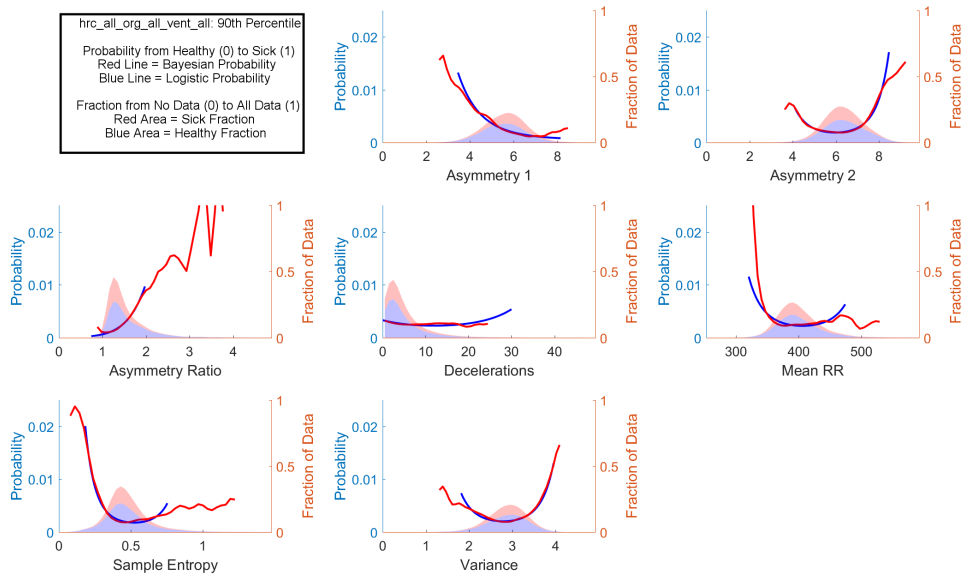


Figure A.17

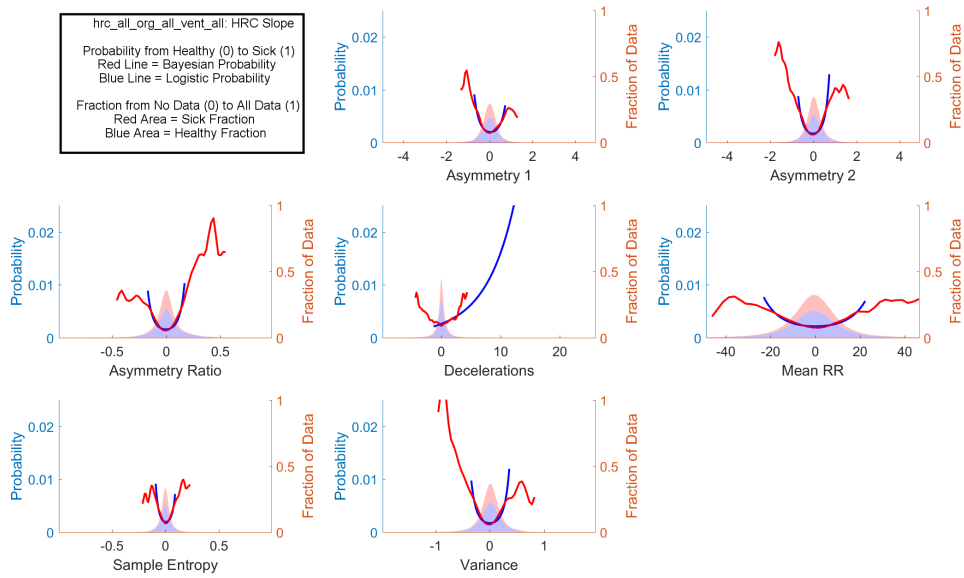
hrc_all_org_all_vent_all: 90th Percentile
 Probability from Healthy (0) to Sick (1)
 Red Line = Bayesian Probability
 Blue Line = Logistic Probability
 Fraction from No Data (0) to All Data (1)
 Red Area = Sick Fraction
 Blue Area = Healthy Fraction



All red and blue areas are additive, i.e., no areas are hidden behind one another.

Figure A.18

hrc_all_org_all_vent_all: HRC Slope
 Probability from Healthy (0) to Sick (1)
 Red Line = Bayesian Probability
 Blue Line = Logistic Probability
 Fraction from No Data (0) to All Data (1)
 Red Area = Sick Fraction
 Blue Area = Healthy Fraction



All red and blue areas are additive, i.e., no areas are hidden behind one another.

Figure A.19

A.7 HeRO Score Figures

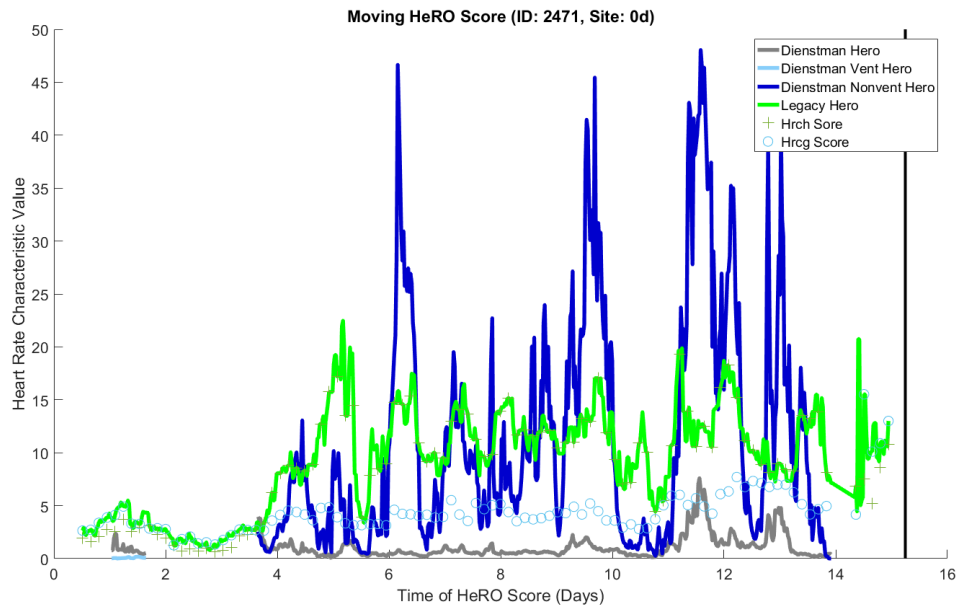


Figure A.20

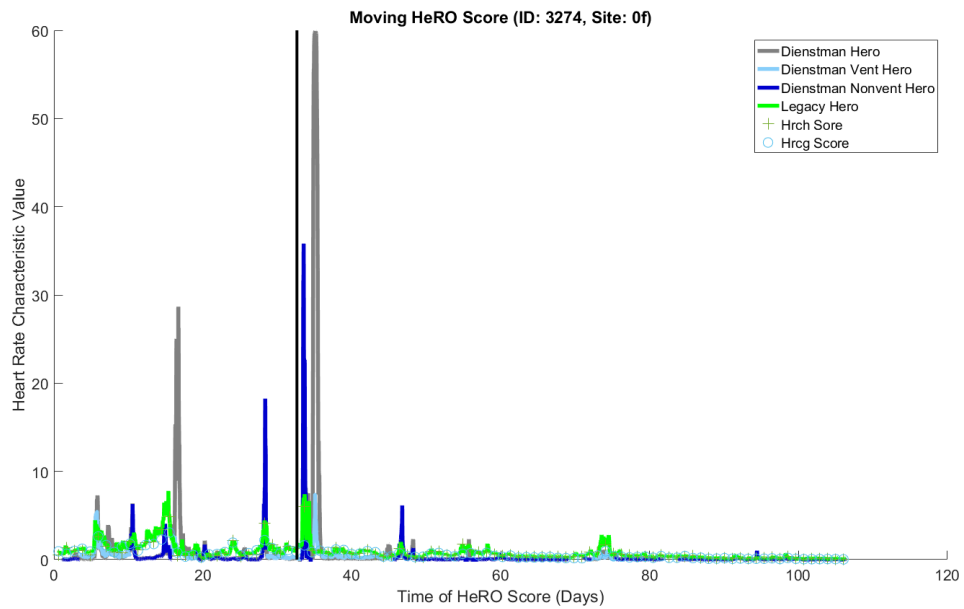


Figure A.21

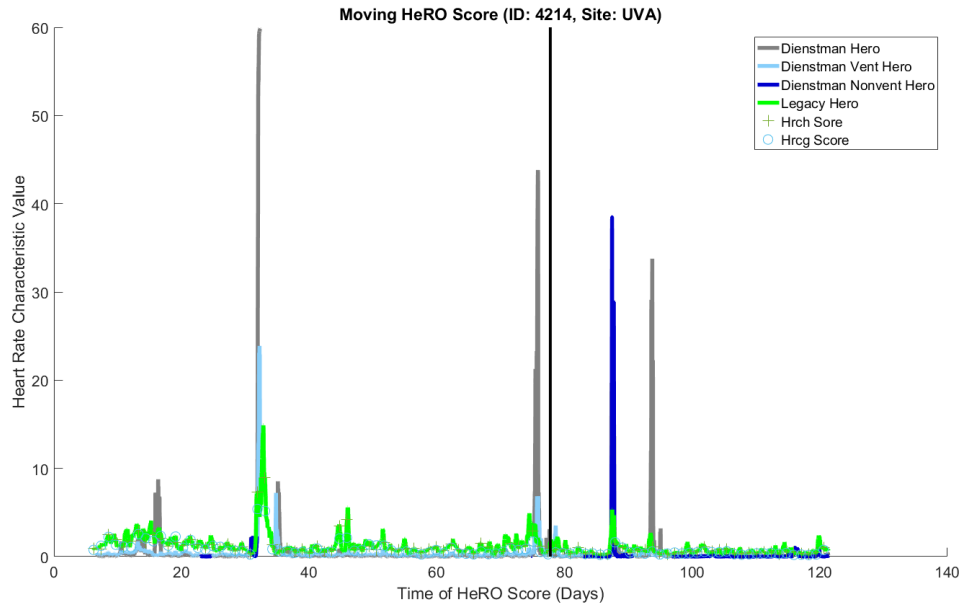


Figure A.22

Appendix B

Matlab Programs

B.1 Result Files

Listing B.1: Dienstman_submit_paralle.txt

```
1 #!/bin/tcsh
2
3 # Check for right number of arguments
4 if ($# != 2) then
5     echo "Usage: Dienstman_submit_parallel begin end"
6     exit 1
7 endif
8
9 if ($1 > $2) then
10    echo "Error: Beginning value is larger than ending value."
11    exit 1
12 endif
13
14 # Generate an array of jobs
15 qsub -N Dienstman_parallel -t $1-$2 <<EOF
16 #!/bin/tcsh
17 #PBS -l nodes=1:c9:ppn=1
18 #PBS -l walltime=180:00:00
19 #PBS -j oe
20 #PBS -q matlab
21 cd /sciclone/home00/eddiest/data10/Dienstman_Files
22 matlab -nojvm -nodisplay -r "multiple_result_files(\$PBS_ARRAYID)" > \$PBS_ARRAYID.out
23 EOF
24
25 exit
```

Listing B.2: multiple_result_files.m

```
1 function multiple_result_files(num)
2 % Author: Evan Dienstman
3 % Last Update: 2/24/2017
4 % Email: eddiestman@email.wm.edu
5 % Note: Feel free to email me with questions! If something doesn't
6 % make sense, it might be because I haven't updated the code yet.
7 %
8 % This program takes one batch number and computes the heart
9 % rate characteristics (HRCs) for each half hour of each ID number in
10 % that batch. Batches consists of 20 ID numbers. This architecture was
11 % created in order to run jobs in parallel most efficiently that also
12 % comply with the rules for W&M's SciClone HPCs. If an ID does not
13 % exist, another function will print "Failed: ID" and continue to the
14 % next ID number. Half hours associated with one ID are saved in a file
15 % corresponding to the site and ID. For more info about what each file
16 % looks like, see the documentation for one_result_file.m. All files
17 % with the same site are saved in a folder corresponding to that site.
18 % If the folder has already been created, the patient files will
19 % automatically be saved there. The information used to calculate the
20 % HRCs comes from the HRC and storm files corresponding to the same
21 % site and ID.
22 %
23 % Arguments:
24 %     1. num - the batch number where the IDs in each batch are defined
25 %           within the code below
26 %
27 % Precondtions:
28 %     1. Make sure "preprocessing" is set correctly in this function.
29 %     2. Make sure the directories in this function are set correctly.
```

```

30 %      3. Make sure the file one_result_file.m is in the working
31 %      directory.
32 %
33 % Returns:
34 %      1. This function returns all files for existing IDs in the batch
35 %      saved to the appropriate folder.
36 %
37 % Change this with the proper preprocessing and site.
38 preprocessing = 'Abby';
39 % preprocessing = 'Doug';
40 %
41 % Here, we define which IDs go into which batch.
42 sites = {};
43 step = 20;
44
45 indices_0d = 2153:step:3287;
46 indices_0f = 2810:step:3981;
47 indices_17 = 2261:step:2512;
48 indices_18 = 2309:step:2486;
49 indices_1a = 2044:step:2637;
50 indices_1b = 2220:step:3078;
51 indices_1e = 2231:step:3315;
52 indices_UVA = 4102:step:7651;
53 start_indices = [indices_0d indices_0f indices_17 indices_18 indices_1a indices_1b indices_1e
    indices_UVA];
54
55 for j = 1:length(indices_0d)
56     sites = [sites '0d']; %#ok<*AGROW>
57 end
58
59 for j = 1:length(indices_0f)
60     sites = [sites '0f'];
61 end
62
63 for j = 1:length(indices_17)
64     sites = [sites '17'];
65 end
66
67 for j = 1:length(indices_18)
68     sites = [sites '18'];
69 end
70
71 for j = 1:length(indices_1a)
72     sites = [sites '1a'];
73 end
74
75 for j = 1:length(indices_1b)
76     sites = [sites '1b'];
77 end
78
79 for j = 1:length(indices_1e)
80     sites = [sites '1e'];
81 end
82
83 for j = 1:length(indices_UVA)
84     sites = [sites 'UVA'];
85 end
86
87 N = length(start_indices);
88
89 % After defining the batches, we call the function that will calculate
90 % the HRCs for each ID number in the batch.
91 if num <= N
92     site = sites{num};
93     start_index = start_indices(num);
94
95     % Now we start looping through the IDs and call one_result_file()
96     % to calculate the HRCs for each ID. The HRB and storm file
97     % contain the information needed to calculate the HRCs.
98     for id = start_index:start_index+step-1
99         hrc_directory = [pwd '/Data/Files/Dienstman_Results_' preprocessing '_PP/' site];
100
101         if ~exist(hrc_directory, 'dir')
102             mkdir(hrc_directory);
103         end
104
105         if strcmp(site, 'UVA')
106             hrb_file = [pwd '/Data/Files/HRB-Files/' site '//UVA_id' num2str(id) '_vch1.hrb'];
107         else
108             hrb_file = [pwd '/Data/Files/HRB-Files/' site '//_id' num2str(id) '_vch1.hrb'];
109         end
110
111         storm_file = [pwd '/Data/Files/Coleman_Results/' site '//storm_results_' site '_id'
            num2str(id) '.mat'];
112         save_file = [hrc_directory '//Dienstman_hrc_results_' site '_' num2str(id) '.mat'];
113         one_result_file(hrb_file, storm_file, save_file, id, site, preprocessing)
114     end
115 end
116
117 end

```

Listing B.3: one_result_file.m

```

1 function one_result_file(hrb_file , storm_file , save_file , id , site , preprocessing)
2 % Author: Evan Dienstman
3 % Last Update: 2/24/2017
4 % Email: eddienstman@email.wm.edu
5 % Note: Feel free to email me with questions! If something doesn't
6 % make sense , it might be because I haven't updated the code yet.
7 %
8 % This function computes the heart rate characteristics (HRCs) for each
9 % half hour of one patient. Patients are identified by their site and
10 % ID number. The half hours are saved in one MATLAB file. Each file is
11 % a struct where each field correspond to one HRC. Each row contains
12 % the HRCs for one half hour. More information about the individual
13 % HRCs can be found throughout the code. Parts of this function were
14 % taken from Abigail Flower and Douglas Lake.
15 %
16 % Arguments:
17 % 1. hrb_file - the name of the HRB file which contains the RR
18 % intervals used for calculating the HRCs. More info about the HRB
19 % files and RR intervals can be found within the code.
20 % 2. storm_file - the name of the storm file which contains the
21 % deceleration HRC. More info on the decelerations can be found
22 % within the code.
23 % 3. save_file - the name of file the struct is saved to where the
24 % name also contains the complete pathway to the file
25 % 4. id - the ID number of the patient
26 % 5. site - the string of the site for use in the save file
27 % 6. preprocessing - a string specifying the type of processing the
28 % raw data goes through. Descriptions of the processing can be
29 % found below.
30 %
31 % Preconditions:
32 % 1. The complete pathway to the HRB, storm, and save file must be
33 % included in respective input variables.
34 % 2. Make sure all the functions this program calls are in the same
35 % directory.
36 %
37 % Returns:
38 % 1. This function returns a MATLAB file containing the HRCs for
39 % one baby.
40 %
41 % We first check multiple exception cases before making the save file.
42 if ~exist(hrb_file , 'file');
43     disp(['Failed: ', num2str(id), ' hrb file does not exist.' char(10)])
44 else
45     % Here, we extract the RR intervals from the HRB file. An RR
46     % interval is the time between each heart beat. RR interval time
47     % is measure in miliseconds. All HRCs are calculated from the
48     % RR intervals. For example, the variance is the variance of the
49     % RR intervals.
50     [rr,rrt,drop,info] = gethrb(hrb_file,inf,1);
51     start_ind_pre = find(rr < 1000);
52
53     if isempty(start_ind_pre)
54         disp(['Failed: ', num2str(id), ' hrb file is empty.' char(10)])
55     else
56
57         if ~exist('rrt', 'var')
58             disp(['Failed: ', num2str(id), ' hrb file has no time variable rrt.' char(10)])
59         else
60             % Here we create a blank struct where we will store the
61             % HRCs. Note that the variable max_num_intervals is used to
62             % preallocate the struct. If there are big jumps in the
63             % times of the file, then the number of intervals
64             % (half hours) containing data will be less than the
65             % maximum number of intervals. The actual number of
66             % intervals will only be equal to the max number of
67             % intervals if the time is fairly continuous (no big jumps).
68             % The way this function deals with jumps in the time is by
69             % taking the start time and end time of the HRB file and
70             % creating enough half hours to fill that entire period of
71             % time. However, since there are usually times with no data,
72             % there are usually rows in the struct left unfilled. This
73             % way is a bit unnecessary but keeps the format coherent
74             % with other files already created and provides an easy way
75             % to preallocate the struct. It also provides a good measure
76             % of how many jumps in the time there are by comparing
77             % intervals with data to the total number of intervals in
78             % the structure.
79             start_ind = start_ind_pre(1);
80             k = start_ind;
81             end_file_time = rrt(end)*24*60;
82             start_file_time = rrt(start_ind)*24*60;
83             length_file = end_file_time-start_file_time;
84             max_num_intervals = ceil(length_file/30);
85             save_variable = ['Dienstman_hrc_results_' site '_' num2str(id)];
86
87             % The purpose of these eval statements is so we can
88             % personalize the name of the variable for each patient.
89             % This way, we can pull up multiple variables from
90             % different files and not get confused as to which variable
91             % belongs to which patient.

```

```

95 | eval([save_variable '(1: num2str(max_num_intervals) ) = struct( 'Start-Time', [],
96 | ' 'Asymmetry_1', [], 'Asymmetry_1-10', [], 'Asymmetry_1-50', [], '
97 | 'Asymmetry_1-90', [], 'Asymmetry_1-Slope', [], '...
98 | 'Asymmetry_2', [], 'Asymmetry_2-10', [], 'Asymmetry_2-50', [], '
99 | 'Asymmetry_2-90', [], 'Asymmetry_2-Slope', [], '...
100 | 'Asymmetry_Ratio', [], 'Asymmetry_Ratio-10', [], 'Asymmetry_Ratio-50', [],
101 | 'Asymmetry_Ratio-90', [], 'Asymmetry_Ratio-Slope', [], '...
102 | 'Decelerations', [], 'Decelerations-10', [], 'Decelerations-50', [], '
103 | 'Decelerations-90', [], 'Decelerations-Slope', [], '...
104 | 'Mean-RR', [], 'Mean-RR-10', [], 'Mean-RR-50', [], 'Mean-RR-90', [], '
105 | 'Mean-RR-Slope', [], '...
106 | 'Sample-Entropy', [], 'Sample-Entropy-10', [], 'Sample-Entropy-50', [], '
107 | 'Sample-Entropy-90', [], 'Sample-Entropy-Slope', [], '...
108 | 'Variance', [], 'Variance-10', [], 'Variance-50', [], 'Variance-90', [],
109 | 'Variance-Slope', [], '...
110 | 'Good-Frac', [], '...
111 | 'Extra-Info', []]);')
112 |
113 | % Now we loop through every half hour and calculate the
114 | % HRCs for that half hour. We use k to index the RR
115 | % interval vector and then find the index closest to
116 | % one half hour from k. This gives us the RR intervals
117 | % within the current half hour. We then use these RR
118 | % intervals to calculate the HRCs for the current half
119 | % hour. Note that each half hour will have a different
120 | % number of intervals for many reasons.
121 | for halfhour = 1:max_num_intervals
122 |
123 |     if k < length(rrt)
124 |         start_time = rrt(k);
125 |         end_time = rrt(k)+(30*(1/60)*(1/24));
126 |
127 |         % Here, there are two types of preprocessing (for
128 |         % the RR intervals) to choose from. The the first
129 |         % one (preprocess) was used by Abigail Flower and
130 |         % the second one was used by Prof. Lake. Abigail's
131 |         % replaces bad intervals with interpolated ones and
132 |         % Lake's simply removes them and concatenates the
133 |         % vector.
134 |         if strcmp(preprocessing, 'Abby')
135 |             end_ind_pre = find(abs(end_time-rrt) == min(abs(end_time-rrt)));
136 |             end_ind = end_ind_pre(1);
137 |             half_hour_indices = k:min(end_ind, length(rrt));
138 |
139 |         else
140 |             half_hour_indices = find(rrt > start_time & rrt <= end_time);
141 |             end_ind = half_hour_indices(end);
142 |         end
143 |
144 |         raw_rr_interval_times = rrt(half_hour_indices);
145 |         raw_rr_intervals = rr(half_hour_indices);
146 |         raw_drop_rr_intervals = drop(half_hour_indices);
147 |         good_frac = 1 - sum(raw_drop_rr_intervals)/length(raw_rr_intervals);
148 |
149 |         if strcmp(preprocessing, 'Abby')
150 |             [processed_rr_intervals, ignore1, ignore2, processed_rr_interval_times] =
151 |             preprocess(raw_rr_intervals, raw_drop_rr_intervals); %%ok<ASGLU>
152 |             processed_rr_interval_times = processed_rr_interval_times ./
153 |             (1000*60*60*24) + start_time;
154 |
155 |         else
156 |             processed_rr_intervals = raw_rr_intervals(raw_drop_rr_intervals == 0);
157 |             processed_rr_interval_times = raw_rr_interval_times(raw_drop_rr_intervals
158 |             == 0);
159 |         end
160 |
161 |         extra_info = struct('info', info, 'raw_rr_intervals', raw_rr_intervals, '
162 |         raw_rr_interval_times', raw_rr_interval_times, 'processed_rr_intervals',
163 |         processed_rr_intervals, 'processed_rr_interval_times',
164 |         processed_rr_interval_times);
165 |
166 |         % This is where the HRC's for good half hours are
167 |         % created. Note note that we have already completed
168 |         % all the preprocessing so the variable
169 |         % processed_drop_rr_interval is all zero indicating
170 |         % nothing will be dropped in the call to
171 |         % calchrcx. For more info about each HRCs
172 |         % calculated in calchrcx(), see the documentation
173 |         % for calchrcx().
174 |         if length(processed_rr_intervals)>300
175 |             cflag = [1,1,1,1,0];
176 |             filter = 1;
177 |             processed_drop_rr_intervals = zeros(1,length(processed_rr_intervals));
178 |             hrc_values = calchrcx(processed_rr_intervals, processed_drop_rr_intervals,
179 |             cflag, filter);
180 |
181 |             % Here we extract the HRCs calculated in
182 |             % calchrcx(). Note that variance and sample
183 |             % asymmetry measurements are recorded on a
184 |             % natural log scale.
185 |             variance = num_check(hrc_values(1));
186 |             sampen = num_check(hrc_values(2));
187 |             asym1 = num_check(hrc_values(3));

```

```

174     asym2 = num_check(hrc_values(4));
175     asym_ratio = asym2/asym1;
176     mean_rr = num_check(mean(processed_rr_intervals));
177
178     % Here we get the deceleration HRC from the
179     % storm file.
180     if exist(storm_file, 'file')
181         load_variable = load(storm_file);
182         storm_results = load_variable.storm_results;
183         height_list = storm_results(halfhour).final_height;
184         decels = 0;
185
186         % The variable r1 is the sample asymmetry
187         % measurement for accelerations. In
188         % symbols, r1 = sum[(beat < median -
189         % median)^2] / total_beats.
190         r1 = exp(hrc_values(3));
191
192         for height = height_list.'
193             % We only include decelerations above a
194             % 5*r1 in height. In other words, we
195             % only include decelerations that are 5
196             % times greater than the average
197             % acceleration.
198             if height > 5 * sqrt(r1)
199                 decels = decels + 1;
200             end
201         end
202
203     else
204         decels = NaN;
205     end
206
207     % The next set of HRCs to calculate are the
208     % slopes and percentiles of the HRCs over a
209     % window in the past for the HRCs already
210     % calculated above. Note that we first need to
211     % enter the current HRCs into the struct so
212     % the current values are used when calculating
213     % the slopes and percentiles.
214     hrc_fields = {'Start_Time', 'Good_Frac', 'Extra_Info', 'Asymmetry_1', '
        Asymmetry_2', 'Asymmetry_Ratio', 'Decelerations', 'Mean_RR', '
        Sample_Entropy', 'Variance'};
215     hrc_vector = {start_time good_frac extra_info asym1 asym2 asym_ratio
        decels mean_rr sampen variance};
216
217     for i = 1:length(hrc_fields)
218         if i == 3
219             eval(['save_variable '(' num2str(halfhour) ')'. hrc_fields{i} '=
                extra_info;'])
220         else
221             eval(['save_variable '(' num2str(halfhour) ')'. hrc_fields{i} '= '
                num2str(hrc_vector{i}) ';''])
222         end
223     end
224
225     eval(['time_vector = [' save_variable '(:).Start_Time;'])
226
227     % Now that we have entered the current HRCs
228     % into the struct, we can now calculate the
229     % slopes and percentiles of the HRCs. Once we
230     % calculate these additional measurements,
231     % we end by adding these measurements into the
232     % struct as yet another HRC.
233     for i = 4:length(hrc_fields)
234         eval(['hrc_fields{i} '_Vector = [' save_variable '(:)'. hrc_fields{i} '
                ';'])
235         slope_end_time = start_time;
236         slope_start_time = slope_end_time - 2;
237         percentile_end_time = start_time;
238         percentile_start_time = percentile_end_time - 0.5;
239
240         eval(['slope_vector = ' hrc_fields{i} '_Vector(time_vector >=
                slope_start_time & time_vector <= slope_end_time);'])
241         slope_time_vector = time_vector(time_vector >= slope_start_time &
                time_vector <= slope_end_time);
242         eval(['percentile_vector = ' hrc_fields{i} '_Vector(time_vector >=
                percentile_start_time & time_vector <= percentile_end_time);'])
243         percentile_time_vector = time_vector(time_vector >=
                percentile_start_time & time_vector <= percentile_end_time); %#ok
                <*NASGU>
244
245         if ~isempty(slope_vector) && ~isempty(percentile_vector)
246             slope = slope_calculator(slope_time_vector.', slope_vector. ');
247             percentile10 = num_check(prctile1(percentile_vector,10));
248             percentile50 = num_check(prctile1(percentile_vector,50));
249             percentile90 = num_check(prctile1(percentile_vector,90));
250
251         else
252             slope = NaN;
253             percentile10 = NaN;
254             percentile50 = NaN;
255             percentile90 = NaN;
256     end

```

```

257
258         eval([save_variable '(' num2str(halfhour) ') .' hrc_fields{i} '_Slope =
259             num2str(slope) ';' ])
260     eval([save_variable '(' num2str(halfhour) ') .' hrc_fields{i} '_10 =
261         num2str(percentile10) ';' ])
262     eval([save_variable '(' num2str(halfhour) ') .' hrc_fields{i} '_50 =
263         num2str(percentile50) ';' ])
264     eval([save_variable '(' num2str(halfhour) ') .' hrc_fields{i} '_90 =
265         num2str(percentile90) ';' ])
266     end
267 else
268     % This is where blank values are created for
269     % bad half hours. Note NaN indicates bad
270     % HRC's as opposed to [] which indicates empty
271     % HRC slots created from preallocating the
272     % struct that we never ended up using.
273     eval([save_variable '(' num2str(halfhour) ') = struct(''Start_Time'',
274         start_time, ' ...
275         ''Asymmetry_1'', NaN, ''Asymmetry_1_10'', NaN, ''Asymmetry_1_50'',
276         NaN, ''Asymmetry_1_90'', NaN, ''Asymmetry_1_Slope'', NaN, ' ...
277         ''Asymmetry_2'', NaN, ''Asymmetry_2_10'', NaN, ''Asymmetry_2_50'',
278         NaN, ''Asymmetry_2_90'', NaN, ''Asymmetry_2_Slope'', NaN, ' ...
279         ''Asymmetry_Ratio'', NaN, ''Asymmetry_Ratio_10'', NaN, ''
280         Asymmetry_Ratio_50'', NaN, ''Asymmetry_Ratio_90'', NaN, ''
281         Asymmetry_Ratio_Slope'', NaN, ' ...
282         ''Decelerations'', NaN, ''Decelerations_10'', NaN, ''
283         Decelerations_50'', NaN, ''Decelerations_90'', NaN, ''
284         Decelerations_Slope'', NaN, ' ...
285         ''Mean_RR'', NaN, ''Mean_RR_10'', NaN, ''Mean_RR_50'', NaN, ''
286         Mean_RR_90'', NaN, ''Mean_RR_Slope'', NaN, ' ...
287         ''Sample_Entropy'', NaN, ''Sample_Entropy_10'', NaN, ''
288         Sample_Entropy_50'', NaN, ''Sample_Entropy_90'', NaN, ''
289         Sample_Entropy_Slope'', NaN, ' ...
290         ''Variance'', NaN, ''Variance_10'', NaN, ''Variance_50'', NaN, ''
291         Variance_90'', NaN, ''Variance_Slope'', NaN, ' ...
292         ''Good_Frac'', NaN, ' ...
293         ''Extra_Info'', extra_info);'])
294     end
295     end
296     k = end.ind + 1;
297     end
298     % Finally, we save the struct using the site and id in
299     % the file name.
300     eval(['save ' save_file ' ' save_variable ';' ])
301 end
302 end
303 end
304
305 function xnew = num-check(x)
306 % This function turns any HRC that is negative, -inf, or inf to NaN
307 % because an HRC should never be any of these values. Occurances like
308 % those may happen due to bad data, so we simply mark them as NaN.
309
310 xnew = x;
311
312 if isempty(xnew)
313     xnew = NaN;
314 elseif xnew < 0 || xnew == inf || xnew == -inf
315     xnew = NaN;
316 end
317 end
318
319 function slope = slope_calculator(x,y)
320 % This function calculates the slope for any HRC. Our method of
321 % calculating the slope is to find a linear fit for the HRC data over
322 % a window in the past and then map the slope of that linear fit as
323 % the HRC slope for the current half hour.
324
325 if length(x) < 48
326     slope = NaN;
327 else
328     X = [ones(length(x),1) x];
329     coeff = X\y;
330     slope = coeff(2);
331 end
332 end

```

B.2 CSV Files

Listing B.4: Dienstman_submit_batch.txt

```
1 #!/bin/tcsh
2 #PBS -l nodes=1:c16:ppn=1
3 #PBS -l walltime=48:00:00
4 #PBS -j oe
5 #PBS -N Dienstman_batch
6 #PBS -q matlab
7 cd /sciclone/home2/eddiest/data10/Dienstman_Files
8 module load matlab/R2016b
9 matlab -nodisplay <multiple_event_figures_caller.m >output_file.out
```

Listing B.5: csv_files.m

```
1 % Author: Evan Dienstman
2 % Last Update: 3/3/2016
3 % Email: eddiestman@email.wm.edu
4 % Note: Feel free to email me with questions! If something doesn't
5 % make sense, it might be because I haven't updated the code yet.
6 %
7 % This script takes all the information from the Dienstman_result
8 % files and organizes them into csv files. Each row in the csv file
9 % corresponds to one half hour from one baby. For the information
10 % that goes into each column, see the header variable in the script
11 % below. The script first creates a CSV file containing all the sick and
12 % healthy half hours from every baby. Note that a healthy half hour is
13 % any half hours that occurs 7 days before an event or 3 days after an
14 % event, and a sick half hour is any half hour that occurs within 12
15 % hours before an event. Additional CSV files are then created for
16 % specific groups. For example, one CSV file may contain only
17 % nonventilated sick half hours from organism 1. The CSV files are
18 % then saved to the appropriate directory. Meta info about the CSV
19 % files are also saved at the end.
20 %
21 % Preconditions:
22 % 1. Make sure the Dienstman_result files are in the appropriate
23 % directory.
24 % 2. Makes sure event_matrix.m, demographic_matrix.m,
25 % vent_matrix.m, csv_master_file.m, csv_splitter.m,
26 % csv_avg_hrcs.m, csv_bin_widths.m, csv_indices.m, and
27 % csv_logistic_coeffs.m are in the current working directory.
28 %
29 % Returns:
30 % 1. This script creates many CSV files corresponding to
31 % different groups of half hours. CSV files are saved to the
32 % Data_Files/Dienstman_CSV_Files directory.
33 % 2. This script also saves meta info about the CSV files
34 % (averages, bin widths, indices, and logistic coefficients)
35 % to the current working directory.
36 %
37 clear
38 clc
39
40 % Change this with the proper preprocessing and site.
41 preprocessing = 'Abby';
42 %preprocessing = 'Doug';
43
44 % Here, we create the header for the CSV file. You can use the header
45 % to see what info goes into the CSV file.
46 header = {'Site' 'ID' 'Half_Hour_Time' 'Event_Time' 'Gender' 'Gestational_Age' 'Good_Frac' '
47 Health_Status' 'Organism' 'Study' 'Ventilated' 'Vent_Switched' 'Vent_Copies' 'Weight'};
48 hrcs = {'Asymmetry_1' 'Asymmetry_2' 'Asymmetry_Ratio' 'Decelerations' 'Mean_RR' 'Sample_Entropy' '
49 Variance'};
50 hrc_types = {'', '_10', '_50', '_90', '_Slope'};
51
52 for i = 1:length(hrcs)
53     for j = 1:length(hrc_types)
54         header = [header, {hrcs{i} hrc_types{j}}]; %#ok<AGROW>
55     end
56 end
57
58 % Now we create the directory where we will save the CSV file.
59 save_dir = [pwd '/Data_Files/Dienstman_CSV_Files_' preprocessing '_PP'];
60
61 if ~exist(save_dir, 'dir')
62     mkdir(save_dir)
63 end
64
65 % Here, we create the master CSV file containing all the sick and
66 % healthy half hours.
67 master_csv_matrix = csv_master_file(preprocessing, save_dir, header, hrcs, hrc_types);
68
69 % Below, we index the master matrix by different qualifiers and write
70 % the various matrices to CSV files. Thus, we will have many CSV
71 % files for different groups of half hours. Creating these smaller
72 % CSV files helps future scripts because we won't have to read in the
73 % master CSV file every time.
```



```

72 csv_splitter(header, save_dir, master_csv_matrix)
73
74 % Finally, we call some functions that create some meta info for the
75 % CSV files that we will use later in other scripts.
76 csv_avg_hrcs(preprocessing)
77 csv_bin_widths(preprocessing)
78 csv_indices(preprocessing)
79 csv_logistic_coeffs(preprocessing, 'vent')
80 csv_logistic_coeffs(preprocessing, 'nonvent')
81 csv_logistic_coeffs(preprocessing, 'all')

```

Listing B.6: csv_splitter.m

```

1 function csv_splitter(header, save_dir, master_csv_matrix)
2 % Author: Evan Dienstman
3 % Last Update: 3/24/2016
4 % Email: eddiendienstman@email.wm.edu
5 % Note: Feel free to email me with questions! If something doesn't
6 % make sense, it might be because I haven't updated the code yet.
7 %
8 % This function takes the master csv matrix which contains
9 % all the healthy and sick half hours and creates multiple CSV files
10 % from that matrix. Each CSV corresponds to one particular group.
11 % For example, one CSV file may contain only sick, unventilated,
12 % organism 3 half hours. Separating half hours in multiple files
13 % helps run future programs faster because it takes a very long
14 % time each time we have to load the master CSV file.
15 %
16 % Arguments:
17 % 1. header - the header of the CSV files
18 % 2. save_dir - the directory to save all the CSV files
19 % 3. master_csv_matrix - the matrix containing all the healthy
20 % and sick half hours that we want to split into many CSV files
21 %
22 % Preconditions:
23 % 1. Make sure the variable master_csv_matrix is a matrix and not
24 % a CSV file.
25 %
26 % Returns:
27 % 1. This function saves many CSV files where each file contains
28 % half hours for one specific category. CSV files are saved to
29 % the directory specified by the variable save_dir.
30 %
31 % First, we define some variables.
32 health_col = find(strcmp(header, 'Health_Status'));
33 org_col = find(strcmp(header, 'Organism'));
34 vent_col = find(strcmp(header, 'Ventilated'));
35 sick_strs = {'healthy', 'sick'};
36
37 % Now we will loop through every category we want to divided the
38 % master CSV matrix into. For every category, we index the appropriate
39 % half hours from the master CSV matrix and save those half hours as
40 % one CSV file.
41 for sick = 0:1
42     sick_str = sick_strs{sick+1};
43     file_name = [save_dir '/hrc_' sick_str '_org-all-vent-all.csv'];
44     file_id = fopen(file_name, 'w');
45     fprintf(file_id, '%s', header{1,1:end-1});
46     fprintf(file_id, '%s\n', header{1,end});
47     fclose(file_id);
48
49     indices = find(master_csv_matrix(:, health_col) == sick);
50     one_csv_matrix = master_csv_matrix(indices, :); %#ok<*FNDSB>
51     dlmwrite(file_name, one_csv_matrix, '-append')
52
53     for ventilated = 0:1
54         file_name = [save_dir '/hrc_' sick_str '_org-all-vent_' num2str(ventilated) '.csv'];
55         file_id = fopen(file_name, 'w');
56         fprintf(file_id, '%s', header{1,1:end-1});
57         fprintf(file_id, '%s\n', header{1,end});
58         fclose(file_id);
59
60         indices = find(master_csv_matrix(:, health_col) == sick & master_csv_matrix(:, vent_col) ==
61             ventilated);
62         one_csv_matrix = master_csv_matrix(indices, :);
63         dlmwrite(file_name, one_csv_matrix, '-append')
64
65         for organism = 1:5
66             file_name = [save_dir '/hrc_' sick_str '_org-' num2str(organism) '_vent-' num2str(
67                 ventilated) '.csv'];
68             file_id = fopen(file_name, 'w');
69             fprintf(file_id, '%s', header{1,1:end-1});
70             fprintf(file_id, '%s\n', header{1,end});
71             fclose(file_id);
72
73             indices = find(master_csv_matrix(:, health_col) == sick & master_csv_matrix(:, org_col)
74                 == organism & master_csv_matrix(:, vent_col) == ventilated);
75             one_csv_matrix = master_csv_matrix(indices, :);
76             dlmwrite(file_name, one_csv_matrix, '-append')
77         end
78     end
79 end

```

```
77 |
78 | end
```

Listing B.7: csv_avg_hrcs.m

```
1 function csv_avg_hrcs(preprocessing)
2 % Author: Evan Dienstman
3 % Last Update: 3/3/2016
4 % Email: eddiendienstman@email.wm.edu
5 % Note: Feel free to email me with questions! If something doesn't
6 % make sense, it might be because I haven't updated the code yet.
7 %
8 % This function takes the average of each column in the master CSV file
9 % hrc_all_org_all_vent_all.csv and saves the info to the current
10 % working directory. Because the master CSV file is relatively large,
11 % it is helpful to calculate the averages once and save the info
12 % instead of having to load the CSV file each time we need the average.
13 %
14 % Arguments:
15 %     1. preprocessing - the preprocessing method used when determining
16 %        which csv file to use
17 %
18 % Preconditions:
19 %     1. Make sure the file hrc_all_org_all_vent_all.csv is in the
20 %        appropriate directory.
21 %
22 % Returns:
23 %     1. This function saves the column averages of the master CSV
24 %        file as a Matlab file called avg_hrc_values to the current
25 %        working directory.
26 %
27 csv_file = [pwd '/Data-Files/Dienstman.CSV-Files-' preprocessing '_PP//hrc_all_org_all_vent_all.'
28             'csv'];
29 raw_data = dlmread(csv_file, ',', 1, 0);
30 avg_hrc_values = nanmean(raw_data); %#ok<*NASGU>
31 save([pwd '//avg_hrc_values-' preprocessing '_PP'], 'avg_hrc_values')
32 end
```

Listing B.8: csv_bin_widths.m

```
1 function csv_bin_widths(preprocessing)
2 % Author: Evan Dienstman
3 % Last Update: 3/24/2016
4 % Email: eddiendienstman@email.wm.edu
5 % Note: Feel free to email me with questions! If something doesn't
6 % make sense, it might be because I haven't updated the code yet.
7 %
8 % This function creates the bin widths used for smoothing and
9 % integrating our probability density functions (PDFs) of each heart
10 % rate characteristic (HRC). We use the Freedman-Diaconis method to get
11 % a bin width for each HRC. The Freedman-Diaconis method is bin_width =
12 % 2*IQR*n^(-1/3), where IQR stands for inter-quartile range and n is
13 % the number of observations. When calculating our bin widths for each
14 % HRC, we use the file hrc_sick_org_3_vent_0.csv. We then use these bin
15 % widths to smooth or integrate a PDF regardless of what CSV file the
16 % PDF came from. Our thought process is that we want to be consistent
17 % with our choice of bin widths when making calculations amongst various
18 % groups. We use the file hrc_sick_org_3_vent_0.csv because it is our
19 % smallest sizeable category. However, if the size of the groups are
20 % relatively close, it might sometimes be better to use a bin widths
21 % calculated from the size of the group we are looking at. In short,
22 % sometimes we will use these bin widths and sometimes we won't. The
23 % same bin widths can also be used for bivariate PDFs. Originally, we
24 % thought we needed to change the IQR (which goes from the 25th
25 % percentile to the 75th percentile) to the 15th percentile and
26 % 85th percentile for bivariate PDFs. Our thought process was that in
27 % two-dimensions, we are smoothing and intergrating over a box and we
28 % need to keep the area of that box proportional to the entire area.
29 % However, the Matlab function ksdensity already handles this issue if
30 % you simply pass the univariate bin widths as arguments.
31 %
32 % Arguments:
33 %     1. preprocessing - the preprocessing method used when determining
34 %        which csv file to use
35 %
36 % Preconditions:
37 %     1. Make sure the directories and file names used in the scripts
38 %        are the right ones for the computer you are using.
39 %     2. Make sure the csv file for sick_her_org_3_vent_0 exists.
40 %
41 % Returns:
42 %     1. This function creates a file called bin_widths.mat that saves
43 %        all the HRC bin widths. Note that the width for decels will be
44 %        1 and is not saved in the file.
45 %
46 % Makes sure the name of the csv file used to make the bin widths
47 % matched the csv file you'll use on your computer.
48 load('hrc_indices')
49 csv_dir = [pwd '/Data-Files/Dienstman.CSV-Files-' preprocessing '_PP'];
```

```

50 csv_file = [csv_dir '/hrc_sick_org-3_vent-0.csv'];
51 csv_matrix = dlmread(csv_file, ',', 1, 0);
52
53 % Here, we determine how many HRCs we have.
54 [~, N] = size(csv_matrix);
55 bin_widths = zeros(1,N);
56
57 % Next, we iterate through every HRC and make the bin width for that
58 % HRC. Note that some HRCs are category info (e.g. id number) and
59 % will never be used for integrating or smoothing.
60 for i = 1:N
61
62     % All deceleration widths are 1.
63     if i >= decelerations_index && i <= decelerations_90_index
64         bin_widths(i) = 1;
65
66     else
67         % Here, we extract the HRC vector from the csv file.
68         csv_vector = csv_matrix(:,i);
69
70         % Using the Freedman–Diaconis method, we now calculate the bin
71         % width using the vector of HRC values for the specific
72         % HRC we are on.
73         n = length(csv_vector);
74         q1 = quantile(csv_vector, .25);
75         q3 = quantile(csv_vector, .75);
76         bin_IQR = q3 - q1;
77         bin_widths(i) = 2*bin_IQR*(n^(-1/3));
78     end
79 end
80
81 % Finally, we save the bin width vector.
82 save([pwd '/bin_widths_' preprocessing '_PP'], 'bin_widths')
83 end

```

Listing B.9: csv_indices.m

```

1 function csv_indices(preprocessing)
2 % Author: Evan Dienstman
3 % Last Update: 3/2/2016
4 % Email: eddienstman@email.wm.edu
5 % Note: Feel free to email me with questions! If something doesn't
6 % make sense, it might be because I haven't updated the code yet.
7 %
8 % This function saves the column index of every HRC in the CSV file.
9 % Saving the column indices helps future functions because we only
10 % have to load the index variables instead of having to load a CSV file
11 % and search the header.
12 %
13 % Arguments:
14 %     1. preprocessing – the preprocessing method used when determining
15 %        which csv file to use
16 %
17 % Preconditions:
18 %     1. Make sure the variable csv_file is the appropriate directory
19 %        for the computer.
20 %
21 % Returns:
22 %     1. This function saves the column index of every HRC in a file
23 %        called hrc_indices.mat. Every index has a unique variable
24 %        name taken from the header of the CSV files.
25 %
26 % First, we load a CSV file that we can use to look up the column
27 % indices.
28 csv_dir = [pwd '/Data_Files/Dienstman_CSV_Files_' preprocessing '_PP'];
29 csv_file = [csv_dir '//hrc_all_org_all_vent_all.csv'];
30
31 [~, col_labels] = xlsread(csv_file, '1:1');
32
33 % Next, we loop through every column and save the index number with
34 % a unique variable name corresponding to the HRC of that column.
35 for i = 1:length(col_labels)
36     save_variable = [lower(col_labels{i}) '_index'];
37     eval([save_variable ' = i;'])
38
39     if i == 1
40         save('hrc_indices', save_variable)
41
42     else
43         save('hrc_indices', save_variable, '-append')
44     end
45 end
46
47 end

```

Listing B.10: csv_logistic_coeffs.m

```

1 function csv_logistic_coeffs(preprocessing, type_str)
2 % Author: Evan Dienstman

```

```

3 % Last Update: 4/12/2017
4 % Email: eddienstman@email.wm.edu
5 % Note: Feel free to email me with questions! If something doesn't
6 % make sense, it might be because I haven't updated the code yet.
7 %
8 % This function creates the coefficients used in the logistic HeRO
9 % score model. We have three sets of coefficients: all coeffs, vent
10 % coeffs, and nonvent coeffs. Since we use a quadratic model, we would
11 % have 666 coefficients from our 35 HRCs. However, we have hand picked
12 % a subset of the HRCs for the HeRO score.
13 %
14 % Arguments:
15 %     1. preprocessing - the preprocessing method used when determining
16 %        which csv file to use
17 %     2. type_str - a string indicating if we are calculating the
18 %        coefficients for vent, nonvent, or all half hours
19 %
20 % Preconditions:
21 %     1. Make sure the file hrc_indices.mat is in the current working
22 %        directory.
23 %
24 % Returns:
25 %     1. This functions saves all the coefficients in the file
26 %        Dienstman_coeffs.mat. The file contains three vectors for the
27 %        three sets of coefficients, three u0 values, and three
28 %        coefficient name cells.
29 %
30 % First, we define some constants and variables used later. Note we are
31 % only picking 12 HRCs.
32 load('hrc_indices.mat')
33 variables = [asymmetry_1.10_index asymmetry_1.slope_index...
34             asymmetry_2.90_index asymmetry_2.slope_index...
35             asymmetry_ratio_50_index asymmetry_ratio_slope_index...
36             decelerations_90_index...
37             sample_entropy_10_index sample_entropy_slope_index...
38             variance_10_index variance_90_index variance_slope_index];
39 model = 'quadratic';
40
41 % Here, load the appropriate CSV file.
42 csv_dir = [pwd '/Data_Files/Dienstman_CSV_Files_' preprocessing 'PP'];
43 csv_file = [csv_dir '//hrc_all_org_all_vent_all.csv'];
44 csv_matrix = dlmread(csv_file, ',', 1, 0);
45
46 % Here, we load the header and define some variables. We use the
47 % smallest file since all the headers are the same.
48 [~, variable_names] = xlsread([csv_dir '//hrc_sick_org_5_vent_0.csv'], '1:1');
49
50 % Now we create the coefficients for nonventilated, ventilated, or all
51 % half hours.
52 if strcmp(type_str, 'vent')
53     csv_matrix = csv_matrix(csv_matrix(:, ventilated_index) == 1,:);
54 elseif strcmp(type_str, 'nonvent')
55     csv_matrix = csv_matrix(csv_matrix(:, ventilated_index) == 0,:);
56 end
57
58 % Here, we calculate some variables we will use later. The order of
59 % these lines are very important because I reuse variable names. I
60 % do this to save memory since the matrices are very large.
61 N = length(variables);
62 sick_half_hours = csv_matrix(csv_matrix(:, health_status_index) == 1,:);
63 [num_sick, ~] = size(sick_half_hours);
64 [num_all, ~] = size(csv_matrix);
65 u0 = num_sick/num_all;
66 healthy_half_hours = csv_matrix(csv_matrix(:, health_status_index) == 0,:);
67 csv_matrix = [sick_half_hours; healthy_half_hours];
68 response_vector = csv_matrix(:, health_status_index);
69 csv_matrix = csv_matrix(:, variables);
70
71 % Before we calculate the probabilities, we remove outliers from the
72 % data. For the Bayesian method, removing outliers will not affect
73 % the results because outliers will have very low probabilities.
74 % However, we want to remove outliers for the logistic probability
75 % because we don't want to over fit the data at the tails. For
76 % decelerations, we define the high_fence as 30 because the outlier
77 % method removes too much data. This procedure is strictly empirical
78 % and needs to be analyzed further.
79 for j = 1:N
80     variable = variables(j);
81     temp_data = csv_matrix(:, j);
82     q1 = quantile(temp_data, .25);
83     q3 = quantile(temp_data, .75);
84     IQR = q3 - q1;
85     low_fence = q1 - 1.5*IQR;
86     high_fence = q3 + 1.5*IQR;
87
88     if variable >= decelerations_index && variable <= decelerations_90_index
89         temp_data(temp_data < low_fence | temp_data > 30) = NaN;
90         response_vector(temp_data < low_fence | temp_data > 30) = NaN;
91     else
92         temp_data(temp_data < low_fence | temp_data > high_fence) = NaN;
93         response_vector(temp_data < low_fence | temp_data > high_fence) = NaN;
94     end
95
96     csv_matrix(:, j) = temp_data;

```

```

97 end
98
99 % Now we calculate the probability of illness using logistic
100 % regression. We use the Matlab function fitglm to calculate the
101 % probability. For more information, please see the Matlab
102 % documentation of this function.
103 fit = fitglm(csv_matrix, response_vector, model, 'distribution', 'binomial', 'VarNames', [
104     variable_names(variables), 'Health_Status']);
105 coeffs = fit.Coefficients.Estimate;
106 coeff_names = fit.CoefficientNames;
107
108 % Lastly, we save all the variables we need.
109 if strcmp(type_str, 'vent')
110     vent_coeffs = coeffs; %#ok<*NASGU>
111     vent_u0 = u0;
112     vent_coeff_names = coeff_names;
113     save('Dienstman_coeffs_vent.mat', 'vent_coeffs', 'vent_u0', 'vent_coeff_names')
114 elseif strcmp(type_str, 'nonvent')
115     nonvent_coeffs = coeffs;
116     nonvent_u0 = u0;
117     nonvent_coeff_names = coeff_names;
118     save('Dienstman_coeffs_nonvent.mat', 'nonvent_coeffs', 'nonvent_u0', 'nonvent_coeff_names')
119 else if strcmp(type_str, 'all')
120     all_coeffs = coeffs;
121     all_u0 = u0;
122     all_coeff_names = coeff_names;
123     save('Dienstman_coeffs_all.mat', 'all_coeffs', 'all_u0', 'all_coeff_names')
124 end
125 end

```

B.3 Time Series Figures

Listing B.11: multiple_event_fiugres_caller.m

```

1 % Author: Evan Dienstman
2 % Date: 3/30/2016
3 % Email: eddienstman@email.wm.edu
4 % Note: Feel free to email me with questions! If something doesn't
5 % make sense, it might be because I haven't updated the code yet.
6 %
7 % This script calls the function multiple_event_figures for various
8 % combinations of input arguments. For each call to the function, the
9 % function creates a figure for the average HRC at each half hour seven
10 % days before an event and three days after. The figure also contains
11 % the average of the moving slope as well as the average 10th, 50th,
12 % and 90th percentiles of the HRCs. If the argument plot_str is 'yes',
13 % the function will also produce the same figure for each individual
14 % event. Only events that match the arguments target.org are used in
15 % the average. See the documentation of multiple_event_figures
16 % for more detail.
17 %
18 % Preconditions:
19 %     1. Make sure the Dienstman files and storm files are
20 %        all in their proper directories.
21 %     2. Make sure the files avg_hrc_values.m, event_matrix.mat,
22 %        vent_matrix.m, one_event_hrc.m, hrc_indices.m,
23 %        one_event_plot.m, and multiple_event_figures.m are in the
24 %        working directory.
25 %
26 % Returns:
27 %     1. The script returns a figure plotting the average of each HRC
28 %        for each half hour. The half hours plotted are ones 7 days
29 %        before an event and 3 days after.
30 %     2. If plot_str is 'yes', the function will return the same figure
31 %        for each individual event.
32 %
33 clear
34 clc
35
36 % Change the preprocessing to the one you want to use.
37 preprocessing = 'Abby';
38 %preprocessing = 'Doug';
39
40 % Now we call multiple_event_figures for different combinations of
41 % paramters.
42 multiple_event_figures(1, 'yes', preprocessing)
43 multiple_event_figures(2, 'yes', preprocessing)
44 multiple_event_figures(3, 'yes', preprocessing)
45 multiple_event_figures(4, 'yes', preprocessing)
46 multiple_event_figures(5, 'yes', preprocessing)
47 multiple_event_figures('all', 'yes', preprocessing)

```

Listing B.12: multiple_event_fiugres.m

```

1 function multiple_event_figures(target_org , plot_str , preprocessing)
2 % Author: Evan Dienstman
3 % Date: 3/30/2017
4 % Email: eddienstman@email.wm.edu
5 % Note: Feel free to email me with questions! If something doesn't
6 % make sense, it might be because I haven't updated the code yet.
7 %
8 % This function creates a figure of the average HRC at each half hour
9 % seven days before an event and three days after. The figure also
10 % contains the average of the moving slope as well as the average 10th,
11 % 50th, and 90th percentiles of the HRCs. If plot_str is 'yes', the
12 % function will also produce the same figure for each individual event.
13 % Only events that match the target_org are used in the average.
14 %
15 % Arguments:
16 % 1. target_org - the organism number the user wants to find the
17 % average of (can be 'ALL' for all organisms)
18 % 2. plot_str - a string indicating if the user want to plot the
19 % individual event plots
20 % 3. preprocessing - the preprocessing method used when determining
21 % which result files to use
22 %
23 % Precondtions:
24 % 1. Make sure the Dienstman files and storm files are
25 % all in their proper directories.
26 % 2. Make sure the files avg_hrc_vales.m, event_matrix.mat,
27 % vent_matrix.m, one_event_hrc.m, hrc_indices.m, and
28 % one_event_plot.m are in the working directory.
29 %
30 % Returns:
31 % 1. The function returns a figure plotting the average of each HRC
32 % for each half hour. The half hours plotted are ones 7 days
33 % before an event and 3 days after.
34 % 2. If plot_str is 'yes', the function will return the same figure
35 % for each individual event as well.
36 %
37 % We first check that the save file for the average figure doesn't
38 % already exist.
39 avg_figure_directory = [pwd '/Figure_Files/Dienstman_Event_Figures.' preprocessing '_PP/Averages'
40 ];
41 if ~exist(avg_figure_directory , 'dir')
42 mkdir(avg_figure_directory)
43 end
44
45 avg_save_file_str = [avg_figure_directory '//Dienstman_figure_org-' num2str(target_org) '.fig'];
46
47 if exist(avg_save_file_str , 'file')
48 disp('Error: A file already exists with the save file name. The program stopped because
49 running the program would overwrite the existing file.')
50 return
51 end
52
53 % Next, we load some files and define some variables.
54 load('event_matrix.mat');
55 load('vent_matrix.mat');
56 site_map_keys = {11, 13, 15, 23, 24, 26, 27, 30};
57 site_map_values = {'UVA', '0d', '0f', '17', '18', '1a', '1b', '1e'};
58 site_map = containers.Map(site_map_keys, site_map_values);
59 load(['avg_hrc_values_' preprocessing '_PP'])
60 field_names = {'Asymmetry_1', 'Asymmetry_2', 'Asymmetry_Ratio', ...
61 'Decelerations', 'Mean_RR', 'Sample_Entropy', 'Variance'};
62 hrc_types = {'', '10', '50', '90', 'Slope'};
63 vent_strs = {'vent', 'nonvent'};
64 reverse_vent_strs = {'nonvent', 'vent'}; %#ok< *NASGU>
65
66 % Next, we preallocate some empty structures which we will use to store
67 % the average HRCs.
68 avg_hrc_struct_vent(1:482) = struct('Ventilated', NaN);
69 avg_hrc_struct_nonvent(1:482) = struct('Ventilated', NaN);
70 count_hrc_struct_vent(1:482) = struct('Ventilated', NaN);
71 count_hrc_struct_nonvent(1:482) = struct('Ventilated', NaN);
72
73 for i = 1:length(field_names)
74 for j = 1:length(vent_strs)
75 for k = 1:length(hrc_types)
76 eval([' avg_hrc_struct_' vent_strs{j} '(:).' field_names{i} hrc_types{k} '] = deal(
77 NaN; ' ])
78 eval([' count_hrc_struct_' vent_strs{j} '(:).' field_names{i} hrc_types{k} '] = deal
79 (0); ' ])
80 end
81 end
82 end
83
84 % Next, we iterate through every event in the event_matrix file.
85 for i = 1:length(event_matrix)
86 id = event_matrix(i,1);
87 site_num = event_matrix(i,2);
88 site = site_map(site_num);
89 gest_age = event_matrix(i,4);
90 event_time = event_matrix(i,7);
91 total_age = floor(event_time/7) + gest_age;

```

```

91     organism = event_matrix(i,6);
92     ventilated = event_matrix(i,11);
93     birth_weight = event_matrix(i,3);
94     baby_info = struct('ID', id, 'Site', site, 'Event_Time', event_time, 'Gest_Age', gest_age, '
    Total_Age', total_age, 'Organism', organism, 'Ventilated', ventilated, 'Birth_Weight',
    birth_weight);
95
96     vent_indices = find(vent_matrix(:,1) == site_num & vent_matrix(:,2) == id); %#ok<NODEF>
97     baby_vent_info = vent_matrix(vent_indices,3:4); %#ok<*FNDSB>
98
99     % Here, we continue only if the event matches our target organism.
100    if strcmp(num2str(target_org), num2str(organism)) || strcmp(target_org, 'all')
101        Dienstman_file = [pwd '/Data_Files/Dienstman_Results_' preprocessing '_PP/' site '//
    Dienstman_hrc_results_' site '_' num2str(id) '.mat'];
102
103        % Here, we stop the entire function if a file is missing.
104        if ~exist(Dienstman_file, 'file')
105            disp(['Failed: ', site, ', ', num2str(id), ' files or directories do not exist.'])
106            continue
107        end
108
109        % Otherwise, we extract the info for this event.
110        one_hrc_struct = one_event_hrc(baby_info, baby_vent_info, Dienstman_file);
111
112        % If plot_str is 'yes', we create a figure for this individual
113        % event.
114        if strcmp(plot_str, 'yes')
115            one_event_figure = one_event_plot(avg_hrc_values, baby_info, one_hrc_struct);
116            figure_directory = [pwd '/Figure_Files/Dienstman_Event_Figures_' preprocessing '_PP/'
    site];
117
118            if ~exist(figure_directory, 'dir')
119                mkdir(figure_directory)
120            end
121
122            save_file_str = [figure_directory '//Dienstman_figure_' site '_' num2str(id) '_'
    num2str(round(event_time)) '.fig'];
123            hgsave(one_event_figure, save_file_str, '-v7.3')
124        end
125
126        % Here, we update the avg_structs. This is very dense code,
127        % so I apologize that it's hard to read.
128        one_nonvent_indices = find([one_hrc_struct(:).Ventilated] == 0);
129        one_vent_indices = find([one_hrc_struct(:).Ventilated] == 1);
130
131        for x = 1:length(field_names)
132
133            for y = 1:length(vent_strs)
134
135                for z = 1:length(hrc_types)
136                    eval(['temp_vector = [one_hrc_struct(:).' field_names{x} hrc_types{z} '];'])
137                    eval(['temp_vector(one_ reverse_vent_strs{y} _indices) = NaN;'])
138                    eval(['temp_vector = num2cell(nansum([avg_hrc_struct_' vent_strs{y} '(:).'
    field_names{x} hrc_types{z} ']; temp_vector));'])
139                    eval(['[avg_hrc_struct_' vent_strs{y} '(:).' field_names{x} hrc_types{z} ']'
    = deal(temp_vector{:});'])
140
141                    eval(['temp_vector = +isnan([one_hrc_struct(:).' field_names{x} hrc_types{z}
    ']);'])
142                    eval(['temp_vector(one_ reverse_vent_strs{y} _indices) = NaN;'])
143                    eval(['temp_vector = num2cell(nansum([count_hrc_struct_' vent_strs{y} '(:).'
    field_names{x} hrc_types{z} ']; temp_vector));'])
144                    eval(['[count_hrc_struct_' vent_strs{y} '(:).' field_names{x} hrc_types{z} ']'
    = deal(temp_vector{:});'])
145                end
146            end
147        end
148    end
149 end
150
151 % Here, we divide by the total number of half hours used for each
152 % index to get the average.
153 for i = 1:length(field_names)
154
155     for j = 1:length(vent_strs)
156
157         for k = 1:length(hrc_types)
158             eval(['temp_avg_vector = num2cell([avg_hrc_struct_' vent_strs{j} '(:).' field_names{i}
    hrc_types{k} ']' ./ [count_hrc_struct_' vent_strs{j} '(:).' field_names{i}
    hrc_types{k} ']);'])
159             eval(['[avg_hrc_struct_' vent_strs{j} '(:).' field_names{i} hrc_types{k} ']' = deal(
    temp_avg_vector{:});'])
160         end
161     end
162 end
163
164 % Lastly, we plot and save the average figure.
165 avg_event_figure = one_event_plot(avg_hrc_values, count_hrc_struct_vent, avg_hrc_struct_vent,
    count_hrc_struct_nonvent, avg_hrc_struct_nonvent, target_org);
166 hgsave(avg_event_figure, avg_save_file_str, '-v7.3')
167 end

```

Listing B.13: one_event_hrc.m

```

1 function hrc_fig_struct = one_event_hrc(baby_info, baby_vent_info, Dienstman_file)
2 % Author: Evan Dienstman
3 % Date: 3/30/2017
4 % Email: eddienstman@email.wm.edu
5 % Note: Feel free to email me with questions! If something doesn't
6 % make sense, it might be because I haven't updated the code yet.
7 %
8 % This function takes a Dienstman_hrc_results file and find the
9 % HRCs in each half hour seven days before and three days after a
10 % septic event. This function also determines if each half hour is
11 % ventilated or not ventilated. All this information is then stored
12 % in a struct.
13 %
14 % Arguments:
15 % 1. baby_info - a struct containing the id, site, event time,
16 % gestational age, total age, organism, and ventilation of the
17 % patient
18 % 2. baby_vent_info - a matrix containing the start and end times
19 % of each period the baby was ventilated
20 % 3. Dienstman_file - the Dienstman_hrc_results file containing the
21 % HRCs for each half hour
22 %
23 % Returns:
24 % 1. hrc_fig_struct - a structure containing the HRCs and
25 % ventilation status of each half hour seven days before the
26 % event and three days after the event.
27 %
28 % First, we define some variables.
29 id = baby_info.ID;
30 site = baby_info.Site;
31 event_time = baby_info.Event_Time;
32 num_half_hours = 482;
33 field_names = {'Asymmetry_1', 'Asymmetry_2', 'Asymmetry_Ratio', ...
34 'Decelerations', 'Mean_RR', 'Sample_Entropy', 'Variance'};
35 hrc_types = {'', '10', '150', '190', 'Slope'};
36 %
37 % Next, we preallocate a struct where we will store all the HRC info
38 % for seven days before the event and 3 days after the event.
39 hrc_fig_struct(1:num_half_hours) = struct('Ventilated', NaN);
40 %
41 for i = 1:length(field_names)
42     for j = 1:length(hrc_types)
43         eval(['[hrc_fig_struct(:).' field_names{i} hrc_types{j} ' ] = deal(NaN); ' ])
44     end
45 end
46 %
47 % Here, we calculate the time window and load in the corresponding
48 % result file.
49 event_window = (event_time-7):(1/48):(event_time+2/48+3);
50 load_variable = load(Dienstman_file); %#ok<NASGU>
51 eval_string = ['hrc_results_struct = load_variable.Dienstman_hrc_results_' site '_' num2str(id) '];
52 eval(eval_string);
53 %
54 % We now loop through every half hour in the time window and record the
55 % HRCs for each half hour in a struct.
56 for i = 1:length(hrc_results_struct)
57     hrc_entry = hrc_results_struct(i);
58 %
59 % If the half hour falls within the time window, we look up the
60 % corresponding half hour index and record the HRCs associated
61 % with that index.
62 if ~isempty(hrc_entry.Start_Time) && event_window(1) < hrc_entry.Start_Time && hrc_entry.
63     Start_Time < event_window(end)
64     index = find(event_window <= hrc_entry.Start_Time, 1, 'last');
65 %
66     for j = 1:length(field_names)
67         for k = 1:length(hrc_types)
68             variable = [field_names{j} hrc_types{k}];
69             eval(['hrc_fig_struct(index).' variable ' = hrc_entry.' variable ';'])
70         end
71     end
72 %
73 % Lastly, we find the ventilation status for the half hour.
74 if ~isempty(find(hrc_entry.Start_Time > baby_vent_info(:,1) & hrc_entry.Start_Time <
75     baby_vent_info(:,2), 1))
76     hrc_fig_struct(index).Ventilated = 1;
77 else
78     hrc_fig_struct(index).Ventilated = 0;
79 end
80 end
end
end

```

Listing B.14: one_event_plot.m

```

1 function [baby_figure] = one_event_plot(avg_hrc_values, baby_info_1, hrc_struct_1, baby_info_2,
2     hrc_struct_2, organism) %#ok<*INUSL>
3 % Author: Evan Dienstman
4 % Date: 3/30/2016

```



```

4 % Email: eddienstman@email.wm.edu
5 % Note: Feel free to email me with questions! If something doesn't
6 % make sense, it might be because I haven't updated the code yet.
7 %
8 % This function creates a nice plot of the HRCs seven days before and
9 % three days after a septic event for. The plot contains seven subplots,
10 % one for each HRC. Each subplot contains five curves, one for each HRC
11 % subcategory. If the figure is for a single event, the color of the
12 % curve changes from dark to light if the baby is ventilated or not
13 % ventilated, respectively. If the figure is for an average, each plot
14 % contains a light and dark curve for each HRC, giving a total of ten
15 % curves per subplot. Each subplot also contains a horizontal line
16 % representing the average HRC for all healthy and sick half hours from
17 % all babies for comparison. Lastly, each subplot contains a vertical
18 % line representing the time of the event. Note that the left y axis is
19 % used for the slope.
20 %
21 % Arguments:
22 % 1. avg_hrc_values - the average value of the seven HRCs for all
23 % sick and healthy half hours
24 % 2. baby_info_1 - a structure containing the id, site, event time,
25 % gestational age, total age, organism, and ventilation of the
26 % patient (for one event) or a structure containing the
27 % ventilated counts (for an average)
28 % 3. hrc_struct_1 - a structure containing the HRCs seven days
29 % before the event and three days after the event (for one event)
30 % or a structure containing the ventilated average HRCs (for an
31 % average)
32 % 4. baby_info_2 - a structure containing the non-ventilated
33 % counts (for an average)
34 % 5. hrc_struct_2 - a structure containing the non-ventilated
35 % average HRCs (for an average)
36 % 6. organism - the number of the organism (for an average)
37 %
38 % Returns:
39 % 1. baby_figure - a Matlab figure containing the seven subplots
40 % for each HRC seven days before and three days after a septic
41 % event
42 %
43 % First, we load the hrc indices and define some variables.
44 load('hrc_indices.mat')
45 field_names = {'Asymmetry_1', 'Asymmetry_2', 'Asymmetry_Ratio', ...
46 'Decelerations', 'Mean_RR', 'Sample_Entropy', 'Variance'};
47 hrc_types = {'', '_10', '_50', '_90', '_Slope'};
48 left_ylim_vector = {[0 10], [0 10], [0 10], [0 20], [250 550], [0 1.5], [0 5]};
49 right_ylim_vector = {[ -1 1], [ -1 1], [ -1 1], [ -1 1], [ -20 20], [ -1 1], [ -1 1]};
50 nonvent_color_vector = {[135/255, 206/255, 250/255], [0/255, 250/255, 154/255], ...
51 [221/255, 160/255, 221/255], [240/255, 230/255, 140/255], [240/255, 128/255, 128/255]};
52 vent_color_vector = {[0, 0, 128/255], [0, 128/255, 0], [148/255, 0, 211/255], ...
53 [255/255, 140/255, 0], [128/255, 0, 0]};
54 window = 5; %%ok< *NASGU>
55 avg_type = 's';
56 index_vector = -336:145;
57
58 % Next, we get some information for the legend in the top left of the
59 % figure. For one event, we get info on the baby. For averages, we
60 % get info on the average number of half hours used.
61 if isfield(baby_info_1, 'ID')
62     id = baby_info_1.ID;
63     site = baby_info_1.Site;
64     event_time = baby_info_1.Event_Time;
65     gest_age = baby_info_1.Gest_Age;
66     total_age = baby_info_1.Total_Age;
67     organism = baby_info_1.Organism;
68     ventilated = baby_info_1.Ventilated;
69     vent = [hrc_struct_1(:).Ventilated];
70
71 else
72     temp_vent_count_1 = 0;
73     temp_vent_count_2 = 0;
74     temp_nonvent_count_1 = 0;
75     temp_nonvent_count_2 = 0;
76
77     % This loop calculates the average number of half hours used when
78     % calculating the average HRC for each HRC index relative to the event.
79     % In short, we add together the half hours used for each index
80     % across all HRCs and then divide by the total number of indices
81     % across all HRCs. This number just gives us a rough understanding
82     % of how many half hours were used for the average so don't worry
83     % if you don't completely understand this part.
84     for i = 1:length(field_names)
85
86         for j = 1:length(hrc_types)
87             variable = [field_names{i} hrc_types{j}];
88             eval(['vent_half_hour_vector = [baby_info_1(:).' variable '];'])
89             temp_vent_count_1 = temp_vent_count_1 + sum(vent_half_hour_vector);
90             temp_vent_count_2 = temp_vent_count_2 + length(vent_half_hour_vector);
91
92             eval(['nonvent_half_hour_vector = [baby_info_2(:).' variable '];'])
93             temp_nonvent_count_1 = temp_nonvent_count_1 + sum(nonvent_half_hour_vector);
94             temp_nonvent_count_2 = temp_nonvent_count_2 + length(nonvent_half_hour_vector);
95         end
96     end
97

```

```

98     vent_mean_half_hours = temp_vent_count_1 / temp_vent_count_2;
99     nonvent_mean_half_hours = temp_nonvent_count_1 / temp_nonvent_count_2;
100 end
101
102 % We now start creating the figure.
103 baby_figure = figure('Position', [50,50,1600,900]);
104 set(baby_figure, 'color', 'w');
105
106 % Here, we loop through all the subplots. Each subplot corresponds to
107 % one HRC and contains five curves for each of the HRC subcategories.
108 % For one event, the curves will change from dark to light to
109 % represent ventilated and nonventilated, respectively. For an
110 % average, we plot ventilated and nonventilated separately for a
111 % total of 10 curves per subplot.
112 for i = 1:length(field_names)
113     subplot(3,3,i+1)
114     hold on
115
116     for j = 1:length(hrc_types)
117         variable = [field_names{i} hrc_types{j}];
118
119         % We extract the vectors for plotting here. Note that we
120         % smooth the vectors for plotting purposes only.
121         if isfield(baby_info_1, 'ID')
122             eval(['vent_plot_vector = tsmovavg([hrc_struct_1(:).' variable '], avg_type, window);'
123                 ])
124             eval(['nonvent_plot_vector = tsmovavg([hrc_struct_1(:).' variable '], avg_type, window'
125                 ');'])
126             vent_plot_vector(vent == 0) = NaN; %ok<*AGROW>
127             nonvent_plot_vector(vent == 1) = NaN;
128
129         else
130             eval(['vent_plot_vector = tsmovavg([hrc_struct_1(:).' variable '], avg_type, window);'
131                 ])
132             eval(['nonvent_plot_vector = tsmovavg([hrc_struct_2(:).' variable '], avg_type, window'
133                 ');'])
134
135         end
136
137         % If we are plotting the slope subcategory, we use a different
138         % y axis.
139         if j ~= 5
140             yyaxis('left')
141             plot(index_vector, vent_plot_vector, '-', 'LineWidth', 2, 'Color', vent_color_vector{j}
142                 )
143             plot(index_vector, nonvent_plot_vector, '-', 'LineWidth', 2, 'Color',
144                 nonvent_color_vector{j})
145
146         % Else, we plot all the other subcategories on the right y axis.
147         else
148             yyaxis('right')
149             plot(index_vector, vent_plot_vector, '-', 'LineWidth', 2, 'Color', vent_color_vector{j}
150                 )
151             plot(index_vector, nonvent_plot_vector, '-', 'LineWidth', 2, 'Color',
152                 nonvent_color_vector{j})
153             ylim(right_ylim_vector{i})
154         end
155     end
156
157     % Finally, we add more info to the subplots including a horizontal
158     % line for the average HRC across all half hours from all babies and
159     % a vertical line for the time of the event.
160     yyaxis('left')
161     ylim(left_ylim_vector{i})
162     hrc_index = lower([field_names{i} '_index']);
163     eval(['plot(index_vector, ones(1,length(index_vector))*avg_hrc_values(' hrc_index '), '--', '
164         'LineWidth', 2, 'Color', [0, 0, 0])'])
165     line([0 0], left_ylim_vector{i}, 'Color', [0,0,0], 'LineWidth', 2);
166     title(regexp(field_names{i}, '-', ''), 'FontSize', 13)
167     xlim([-336 145])
168     set(gca, 'FontSize', 14)
169     hold off
170 end
171
172 % Finally, we create the legend at the top left of the figure.
173 title_frame = uicontrol('style', 'frame');
174 set(title_frame, 'Position', [208, 635, 320, 195], 'BackgroundColor', [0 0 0])
175
176 if isfield(baby_info_1, 'ID')
177     round_time = round(event_time);
178     baby_title_string = [
179         ' Patient ID Number: ' num2str(id) char(10)...
180         ' Site Code: ' site char(10)...
181         ' Days of Age: ' num2str(round_time) char(10)...
182         ' Gestational Age (Weeks): ' num2str(gest_age) char(10)...
183         ' Total Age (Weeks): ' num2str(total_age) char(10)...
184         ' Organism Number: ' num2str(organism) char(10)...
185         ' Ventilation at Event: ' num2str(ventilated) char(10) char(10)...
186         ' Blue: Raw, Green: 10th, Purple: 50th, ' char(10)...
187         ' Orange: 90th, Red: Slope ' char(10) char(10)...
188         ' Dark: Veniltaed, Ligth: Nonventilaed'];
189 else
190     baby_title_string = [
191         ' Organism Number: ' num2str(organism) char(10) char(10) char
192         (10)...
193         ' Mean Num of Vent Half Hours Used: ' num2str(vent_mean_half_hours) char(10)...

```

```

182     'Mean Num of Nonvent Half Hours Used: ' num2str(nonvent_mean_half_hours) char(10) char(10)
183     ' ...
184     ' Blue: Raw, Green: 10th, Purple: 50th,' char(10)...
185     ' Orange: 90th, Red: Slope' char(10) char(10)...
186     ' Dark: Veniltaed, Ligth: Nonventilaed'];
187 end
188 baby_title = uicontrol('style','text');
189 set(baby_title, 'String', baby_title_string, 'Position', [211, 638, 314, 189], 'FontSize', 9, '
    BackgroundColor', [1 1 1])
190
191 axes_note_string = 'All x-axes show the half hour index. Thus, there are 482 half hours
    representing the 10 day window.';
192 axes_note = uicontrol('style','text');
193 set(axes_note, 'String', axes_note_string, 'Position', [520,-20,600, 50], 'FontSize', 10, '
    BackgroundColor', [1 1 1])
194 end

```

B.4 Univariate PDF Figures

Listing B.15: multiple_univariate_pdf_figures.m

```

1 % Author: Evan Dienstman
2 % Last Update: 3/23/2017
3 % Email: eddienstman@email.wm.edu
4 % Note: Feel free to email me with questions! If something doesn't
5 % make sense, it might be because I haven't updated the code yet.
6 %
7 % This script creates figures for the univariate PDFs of HRCs from
8 % various categories. For more information about the figures, please
9 % see the documentation for one_univariate_pdf_figure.m.
10 %
11 % Preconditions:
12 % 1. Make sure the files in file_list have the correct name.
13 % 2. Make sure the color_list, color_name_list, and file_list
14 % used in one call to one_univariate_pdf_figure have the same
15 % length.
16 % 3. Make sure the file one_univariate_pdf_figure.m is in the
17 % current working directory.
18 %
19 % Returns:
20 % 1. This script will create five figures for each call to
21 % one_bivariate_pdf_figure. For info about where the files are
22 % saved, see the documentation for one_bivariate_pdf_figure.
23 %
24 clear
25 clc
26
27 % Change the preprocessing to the one you want to use.
28 preprocessing = 'Abby';
29 % preprocessing = 'Doug';
30
31 % First, we define the colors of the PDFs in our figure.
32 color_list_1 = {[0, 0, 128/255], [135/255, 206/255, 250/255], [128/255, 0, 0], [240/255, 128/255,
    128/255]};
33 color_list_2 = {[0, 0, 1], [1, 0, 0], [0, 100/255, 0], [128/255, 0, 128/255], [255/255, 165/255,
    0], [0, 0, 0]};
34 color_list_3 = {[0, 0, 128/255], [128/255, 0, 0]};
35
36 color_name_list_1 = {'Dark Blue', 'Light Blue', 'Dark Red', 'Light Red'};
37 color_name_list_2 = {' Blue', ' Red', ' Green', ' Purple', ' Orange', ' Black'};
38 color_name_list_3 = {'Dark Blue', 'Dark Red'};
39
40 % Next, we create lists of CSV files that we want to plot to the
41 % same figure.
42 file_list_1 = {'hrc_healthy_org-1_vent-0.csv', 'hrc_healthy_org-1_vent-1.csv', '
    hrc_sick_org-1_vent-0.csv', 'hrc_sick_org-1_vent-1.csv'};
43 file_list_2 = {'hrc_healthy_org-2_vent-0.csv', 'hrc_healthy_org-2_vent-1.csv', '
    hrc_sick_org-2_vent-0.csv', 'hrc_sick_org-2_vent-1.csv'};
44 file_list_3 = {'hrc_healthy_org-3_vent-0.csv', 'hrc_healthy_org-3_vent-1.csv', '
    hrc_sick_org-3_vent-0.csv', 'hrc_sick_org-3_vent-1.csv'};
45 file_list_4 = {'hrc_healthy_org-4_vent-0.csv', 'hrc_healthy_org-4_vent-1.csv', '
    hrc_sick_org-4_vent-0.csv', 'hrc_sick_org-4_vent-1.csv'};
46 file_list_5 = {'hrc_healthy_org-5_vent-0.csv', 'hrc_healthy_org-5_vent-1.csv', '
    hrc_sick_org-5_vent-0.csv', 'hrc_sick_org-5_vent-1.csv'};
47 file_list_6 = {'hrc_sick_org-1_vent-0.csv', 'hrc_sick_org-2_vent-0.csv', 'hrc_sick_org-3_vent-0.
    csv', 'hrc_sick_org-4_vent-0.csv', 'hrc_sick_org-5_vent-0.csv', 'hrc_all_org-all_vent-all.csv'
    };
48 file_list_7 = {'hrc_healthy_org-all_vent-0.csv', 'hrc_sick_org-all_vent-0.csv'};
49 file_list_8 = {'hrc_healthy_org-all_vent-1.csv', 'hrc_sick_org-all_vent-1.csv'};
50 file_list_9 = {'hrc_healthy_org-all_vent-all.csv', 'hrc_sick_org-all_vent-all.csv'};
51
52 % Finally, we call one_univariate_pdf_figure for each file_list.
53 % Every call to one_univariate_pdf_figure will create five figures
54 % (one for each HRC subcategory). Each of the five figures will
55 % contain PDFs from the files in file_list.

```

```

56 one_univariate_pdf_figure('hrc.all.org-1.vent.all', file_list_1, color_list_1, color_name_list_1,
    preprocessing)
57 one_univariate_pdf_figure('hrc.all.org-2.vent.all', file_list_2, color_list_1, color_name_list_1,
    preprocessing)
58 one_univariate_pdf_figure('hrc.all.org-3.vent.all', file_list_3, color_list_1, color_name_list_1,
    preprocessing)
59 one_univariate_pdf_figure('hrc.all.org-4.vent.all', file_list_4, color_list_1, color_name_list_1,
    preprocessing)
60 one_univariate_pdf_figure('hrc.all.org-5.vent.all', file_list_5, color_list_1, color_name_list_1,
    preprocessing)
61 one_univariate_pdf_figure('hrc.sick.org-all.vent.0', file_list_6, color_list_2, color_name_list_2,
    preprocessing)
62 one_univariate_pdf_figure('hrc.all.org-all.vent.0', file_list_7, color_list_3, color_name_list_3,
    preprocessing)
63 one_univariate_pdf_figure('hrc.all.org-all.vent.1', file_list_8, color_list_3, color_name_list_3,
    preprocessing)
64 one_univariate_pdf_figure('hrc.all.org-all.vent.all', file_list_9, color_list_3, color_name_list_3,
    preprocessing)

```

Listing B.16: one.univariate_pdf_figure.m

```

1 function one_univariate_pdf_figure(save_file_str, file_list, color_list, color_name_list,
    preprocessing)
2 % Author: Evan Dienstman
3 % Last Update: 3/23/2017
4 % Email: eddiendienstman@email.wm.edu
5 % Note: Feel free to email me with questions! If something doesn't
6 % make sense, it might be because I haven't updated the code yet.
7 %
8 % This function creates univariate probability density functions (PDFs)
9 % for each heart rate characteristics (HRC) of each csv file given in
10 % the variable file_list. The PDFs are then all plotted onto one
11 % of five figures. Each of the five figures contains the PDFs for one
12 % HRC subtype (raw, 10th, 50th, 90th, and slope). One each figure,
13 % there are seven subplots that corresponds to the seven HRCs. Each
14 % subplot contains a PDF for every CSV in the variable file_list.
15 % There is also a text box at the top of the figure that acts as a
16 % legend for each subplot. This text box also contains information
17 % on how many babies and half hours were used for each PDF. Note that
18 % for any sick category, we take any half hour within 24 half hours
19 % of an event. For any healthy category, we take any half hour
20 % 7 days before an event or 3 days after event. If a baby never had an
21 % event, then all half hours are healthy. The PDFs were created using
22 % kernel density estimation. For more details about this method, see
23 % the comments in the code. Each row in the csv files corresponds to
24 % one half hour. Each column contains an HRC or info about the half
25 % hour. Different csv files contain half hours that fall into different
26 % categories from the information in the category columns. If the user
27 % wants to graph the PDFs from different categories of half hours, they
28 % can do so by adding, removing, or changing the corresponding csv
29 % files in the variable file_list.
30 %
31 % Arguments:
32 % 1. save_file_str - the name of the save files
33 % 2. file_list - the CSV files used to make the PDFs
34 % 3. color_list - the RGB colors of the PDFs
35 % 4. color_name_list - the color names of the PDFs
36 % 5. preprocessing - a string indicating which preprocessing method
37 % to use
38 %
39 % Preconditions:
40 % 1. Make sure the directories and file names used in the scripts
41 % are the right ones for the computer you are using.
42 % 2. Make sure file_list, color_list, and color_name_list contain
43 % the same number of elements.
44 % 3. This script will not overwrite any existing files with the
45 % same name. Change the variable save_file to a name that does
46 % not already exist or delete the existing file with the same
47 % name before running this script.
48 %
49 % Returns:
50 % 1. This function creates five figures saved with the name given
51 % by the variable save_file and with the appropriate subcategory
52 % string at the end. The figures are saved to the directory
53 % given by the variable save_dir below.
54 %
55 % Change the directory names to match the directories on your
56 % computer.
57 csv_dir = [pwd '/Data_Files/Dienstman_CSV_Files_' preprocessing '_PP'];
58 save_dir = [pwd '/Figure_Files/Dienstman_Univariate_PDFs_' preprocessing '_PP'];
59 save_file = [save_dir '/' save_file_str '.png'];
60
61 % If the save_dir doesn't already exist, we make the save_dir here.
62 if ~exist(save_dir, 'dir')
63     mkdir(save_dir)
64 end
65
66 % If the save file already exists, we stop the program so we don't
67 % overwrite the file.
68 if exist(save_file, 'file')
69     disp('Error: A file already exists with the save file name. The program stopped because
    running the program would overwrite the existing file.')

```

```

70     return
71 end
72
73 % Here, we define the edges used for plotting.
74 N = 1e3;
75 asym1_edges = linspace(0,10,N);
76 asym2_edges = linspace(0,10,N);
77 asym_ratio_edges = linspace(0,5,N);
78 decel_edges = linspace(0,50,N);
79 mean_rr_edges = linspace(250,550,N);
80 sampen_edges = linspace(0,1.25,N);
81 variance_edges = linspace(0,5,N);
82 edge_vector = {asym1_edges, asym2_edges, asym_ratio_edges, decel_edges, mean_rr_edges,
    sampen_edges, variance_edges};
83
84 asym1_edges_slope = linspace(-5,5,N);
85 asym2_edges_slope = linspace(-5,5,N);
86 asym_ratio_edges_slope = linspace(-1,1,N);
87 decel_edges_slope = linspace(-10,10,N);
88 mean_rr_edges_slope = linspace(-50,50,N);
89 sampen_edges_slope = linspace(-1,1,N);
90 variance_edges_slope = linspace(-2,2,N);
91 edge_vector_slope = {asym1_edges_slope, asym2_edges_slope, asym_ratio_edges_slope,
    decel_edges_slope, mean_rr_edges_slope, sampen_edges_slope, variance_edges_slope};
92
93 % These vectors are used to create the information in the boxes at
94 % the top of the figures.
95 baby_count_vector = zeros(1,length(file_list));
96 mean_half_hour_vector = zeros(1,length(file_list));
97
98 % These cells are used to create the HRC names.
99 hrc_names = {'Asymmetry-1', 'Asymmetry-2', 'Asymmetry-Ratio', 'Decelerations', 'Mean-RR', '
    Sample-Entropy', 'Variance'};
100 hrc_types = {'', '_10', '_50', '_90', '_Slope'};
101
102 % Here, we load the header. We use the smallest file since all the
103 % headers are the same.
104 [~, header_names] = xlsread([csv_dir '//hrc_sick_org_5_vent_0.csv'], '1:1');
105
106 % Next, we loop through every file in file_list and make the PDFs for
107 % each HRC of that file.
108 for i = 1:length(file_list)
109     csv_file = [csv_dir '/' file_list{i}];
110     raw_data = dlmread(csv_file, ',', 1, 0);
111
112     % Next, we calculate the info that goes into the box at the top
113     % of the figure for this csv file.
114     site_index = find(strcmp(header_names, 'Site'));
115     id_index = find(strcmp(header_names, 'ID'));
116     unique_matrix = unique(raw_data(:, [site_index id_index]), 'rows');
117     baby_count_vector(i) = length(unique_matrix);
118     mean_half_hour_vector(i) = length(raw_data)/length(unique_matrix);
119
120     % Now we can loop through every HRC and make the PDF for that
121     % HRC with the half hours of the current CSV file.
122     for j = 1:length(hrc_names)
123
124         for k = 1:length(hrc_types)
125             variable_name = [hrc_names{j} hrc_types{k}];
126
127             if k == 5
128                 edge = edge_vector_slope{j}; % #ok< *NASGU>
129             else
130                 edge = edge_vector{j};
131             end
132
133             % Next, we extract the HRC information from the csv file.
134             column_data = raw_data(:, strcmp(header_names, variable_name));
135
136             % Here, we create the PDFs for each HRC of this file.
137             % We create PDFs using the Matlab function ksdensity,
138             % which uses the kernel density estimation method to
139             % create the PDFs. For the kernel density estimation,
140             % we use an Epanechnikov kernel. We also allow the
141             % function to calculate the optimal bandwidth except for
142             % decelerations where we use a width of 1. At first,
143             % we used the width the Freeman-Diaconis method produced
144             % from the organism 3 file so all the PDFs would have a
145             % consistent width. However, we believe this isn't a
146             % problem if all our n's are relatively close. The
147             % input for ksdensity is a vector of one specific type of
148             % HRC values corresponding to every half hour in the
149             % csv file. Note that since we are creating PDFs, the
150             % area underneath each curve is 1. Consequently, all
151             % the PDFs are normalized and values of the PDF can be
152             % greater than 1. For more information about the kernel
153             % density estimation, Epanechnikov kernel function, and
154             % PDFs, see their respective wikipedia pages.
155             if strcmp(hrc_names{j}, 'Decelerations') || k == 5
156                 eval([variable_name '_Prob(:,i) = ksdensity(column_data, edge, ''kernel'', ''
                    epanechnikov'');'])
157             else
158                 eval([variable_name '_Prob(:,i) = ksdensity(column_data, edge, ''kernel'', ''
                    epanechnikov'', ''width'', 1);'])

```

```

159         end
160     end
161 end
162 end
163
164 % Here, we define some variables used for plotting.
165 x_axis_list = {[0 10], [0 10], [0 3], [0 10], [250 550], [0 1], [0 5]};
166 y_axis_list = {[0 0.75], [0 0.75], [0 5], [0 0.4], [0 0.02], [0 6], [0 1.5]};
167 x_axis_list_slope = {[ -2 2], [ -2 2], [ -0.75 0.75], [ -5 5], [ -50 50], [ -0.5 0.5], [ -1 1]};
168 y_axis_list_slope = {[0 2], [0 2], [0 7.5], [0 1], [0 0.1], [0 15], [0 4]};
169 title_list = {'Raw HRC ', '10th Percentile ', '50th Percentile ', '90th Percentile ', 'HRC Slope '
};
170
171 % Finally, we plot the PDFs. We will create five figures. Each figure
172 % contains the PDFs of one subcategory. The subcategories are raw HRC,
173 % 10th percentile, 50th percentile, 90th percentile, and slope. Each
174 % figure contains seven subplots that corresponds to one specific HRC.
175 % On each subplot, there is a PDF for each CSV file in the variable
176 % file_list.
177 for i = 1:length(hrc_types)
178     probability_figure = figure('Position', [50,50,1600,900]);
179     set(probability_figure, 'color', 'w');
180
181     for j = 1:length(hrc_names)
182         subplot(3,3,j+1)
183         hold on
184
185         if i == 5
186             x = edge_vector_slope{j}.';
187         else
188             x = edge_vector{j}.';
189         end
190
191         variable_name = [hrc_names{j} hrc_types{i}];
192
193         % Here, we plot the PDFs.
194         for k = 1:length(file_list)
195             eval(['y = ' variable_name '_Prob(:,k);'])
196             plot(x, y, 'LineWidth', 2, 'Color', color_list{k})
197         end
198
199         if i == 5
200             xlim(x_axis_list_slope{j})
201             ylim(y_axis_list_slope{j})
202         else
203             xlim(x_axis_list{j})
204             ylim(y_axis_list{j})
205         end
206
207         xlabel(regexprep(hrc_names{j}, '_', ' '), 'FontSize', 16)
208         ylabel('Probability', 'FontSize', 16)
209         set(gca, 'FontSize', 16)
210         hold off
211     end
212
213     % Here, we plot the box at the top of the figure.
214     if length(file_list) == 6
215         baby_title_string = [' Univariate Probability Densities: ' title_list{i}];
216     else
217         baby_title_string = [' Univariate Probability Densities: ' title_list{i} char(10)];
218     end
219
220     for n = 1:length(file_list)
221         file_str = file_list{n};
222         baby_title_string = [baby_title_string char(10) ' ' color_name_list{n} ' = '
file_str(1:end-4) char(10) ' Baby Count = ' num2str(baby_count_vector(n)) ', Mean
Half Hours = ' num2str(mean_half_hour_vector(n), 4)]; %#ok<AGROW>
223     end
224
225     title_frame = uicontrol('style', 'frame');
226     set(title_frame, 'Position', [208, 635, 345, 195], 'BackgroundColor', [0 0 0])
227     baby_title = uicontrol('style', 'text');
228     set(baby_title, 'String', baby_title_string, 'Position', [211, 638, 339, 189], 'FontSize', 8,
'BackgroundColor', [1 1 1])
229
230     % Lastly, we save the figure.
231     print(probability_figure, [save_file(1:end-4) lower(hrc_types{i}) '.png'], '-dpng')
232 end
233 end

```

B.5 Univariate Risk Figures

Listing B.17: multiple_univariate_risk_figures.m

```

1 % Author: Evan Dienstman
2 % Last Update: 3/31/2017
3 % Email: eddienstman@email.wm.edu

```

```

4 % Note: Feel free to email me with questions! If something doesn't
5 % make sense, it might be because I haven't updated the code yet.
6 %
7 % This script creates figures for the univariate risks of HRCs from
8 % various categories. For more information about the figures, please
9 % see the documentation for one_univariate_risk_figure.m.
10 %
11 % Preconditions:
12 %     1. Make sure the files in file_list have the correct names.
13 %     2. Make sure the file one_univariate_risk_figure.m is in the
14 %        current working directory.
15 %     3. Make sure to select the appropriate preprocessing below.
16 %
17 % Returns:
18 %     1. This script will create one figure for each call to
19 %        one_bivariate_pdf_figure. For info about where the files are
20 %        saved, see the documentation for one_univariate_risk_figure.
21 %
22 clear
23 clc
24
25 % Change the preprocessing to the one you want to use.
26 preprocessing = 'Abby';
27 % preprocessing = 'Doug';
28
29 % First, we define different lists of CSV files.
30 file_list_0 = {'hrc_sick_org_all_vent_all.csv', 'hrc_healthy_org_all_vent_all.csv'};
31 file_list_1 = {'hrc_sick_org_all_vent_0.csv', 'hrc_healthy_org_all_vent_0.csv'};
32 file_list_2 = {'hrc_sick_org_all_vent_1.csv', 'hrc_healthy_org_all_vent_1.csv'};
33 file_list_3 = {'hrc_sick_org_1_vent_0.csv', 'hrc_sick_org_2_vent_0.csv', 'hrc_sick_org_ALL_vent_0.
34               .csv'};
34 file_list_4 = {'hrc_sick_org_1_vent_0.csv', 'hrc_sick_org_3_vent_0.csv', 'hrc_sick_org_ALL_vent_0.
35               .csv'};
35 file_list_5 = {'hrc_sick_org_1_vent_0.csv', 'hrc_sick_org_4_vent_0.csv', 'hrc_sick_org_ALL_vent_0.
36               .csv'};
36 file_list_6 = {'hrc_sick_org_1_vent_0.csv', 'hrc_sick_org_5_vent_0.csv', 'hrc_sick_org_ALL_vent_0.
37               .csv'};
37 file_list_7 = {'hrc_sick_org_2_vent_0.csv', 'hrc_sick_org_3_vent_0.csv', 'hrc_sick_org_ALL_vent_0.
38               .csv'};
38 file_list_8 = {'hrc_sick_org_2_vent_0.csv', 'hrc_sick_org_4_vent_0.csv', 'hrc_sick_org_ALL_vent_0.
39               .csv'};
39 file_list_9 = {'hrc_sick_org_2_vent_0.csv', 'hrc_sick_org_5_vent_0.csv', 'hrc_sick_org_ALL_vent_0.
40               .csv'};
40 file_list_10 = {'hrc_sick_org_3_vent_0.csv', 'hrc_sick_org_4_vent_0.csv', 'hrc_sick_org_ALL_vent_0
41                .csv'};
41 file_list_11 = {'hrc_sick_org_3_vent_0.csv', 'hrc_sick_org_5_vent_0.csv', 'hrc_sick_org_ALL_vent_0
42                .csv'};
42 file_list_12 = {'hrc_sick_org_4_vent_0.csv', 'hrc_sick_org_5_vent_0.csv', 'hrc_sick_org_ALL_vent_0
43                .csv'};
43
44 % Next, we call one_univariate_risk_figure for each file_list above.
45 one_univariate_risk_figure(file_list_0, preprocessing)
46 one_univariate_risk_figure(file_list_1, preprocessing)
47 one_univariate_risk_figure(file_list_2, preprocessing)
48 one_univariate_risk_figure(file_list_3, preprocessing)
49 one_univariate_risk_figure(file_list_4, preprocessing)
50 one_univariate_risk_figure(file_list_5, preprocessing)
51 one_univariate_risk_figure(file_list_6, preprocessing)
52 one_univariate_risk_figure(file_list_7, preprocessing)
53 one_univariate_risk_figure(file_list_8, preprocessing)
54 one_univariate_risk_figure(file_list_9, preprocessing)
55 one_univariate_risk_figure(file_list_10, preprocessing)
56 one_univariate_risk_figure(file_list_11, preprocessing)
57 one_univariate_risk_figure(file_list_12, preprocessing)

```

Listing B.18: one_univariate_risk_figure.m

```

1 function one_univariate_risk_figure(file_list, preprocessing)
2 % Author: Evan Dienstman
3 % Last Update: 3/31/2017
4 % Email: eddienstman@email.wm.edu
5 % Note: Feel free to email me with questions! If something doesn't
6 % make sense, it might be because I haven't updated the code yet.
7 %
8 % This script creates univariate risk figures for each heart rate
9 % characteristics (HRC) using the CSV files in file_list. If there are
10 % two files in file_list, then the risk is defined as risk of file_1
11 % relative to file_2. If there are three files in file_list, then the
12 % risk is defined as file_1 relative to file_3 divided by file_2
13 % relative to file_3. For a further explanation of how we define risk,
14 % see the documentation of one_risk_matrix. The risks are then all
15 % plotted onto one figure. There are seven subplots in the figure,
16 % one for each HRC category, and each subplot has a curve that
17 % corresponds to one of the risks for that HRC category. For example,
18 % one subplot will have the risk curves for variance, variance_slope,
19 % variance_l0, etc. There is also a text box at the top of the figure
20 % that acts as a legend for each subplot. This text box also contains
21 % information on how many babies and half hours were used when
22 % calculating the risks. Note that for any sick category, we take any
23 % half hour within 24 half hours of an event. For any healthy category,
24 % we take any half hour 7 days before an event or 3 days after event.
25 % If a baby never had an event, then all half hours are healthy.

```

```

26 %
27 % Arguments:
28 % 1. file_list - the list of csv_files for calculating the risk
29 % 2. preprocessing - a string indicating which preprocessing method
30 % to use
31 %
32 % Preconditions:
33 % 1. Make sure the directories and file names used in the scripts
34 % are the right ones for the computer you are using.
35 % 2. Make sure the files one_risk_matrix.m, one_porb_matrix.m,
36 % and bin_width_maker.m are in the working directory.
37 % 3. The variable file_list can only contain 2 or 3 files.
38 % 4. This script will not overwrite any existing files with the
39 % same name. Change the variable save_file to a name that does
40 % not already exist, or delete the existing file with the same
41 % name before running this script.
42 %
43 % Returns:
44 % 1. This function returns a figure saved with the name given by
45 % the variable save_file that shows the univariate risks
46 % relative to the files in variable file_list for each HRC.
47
48 % Change the directory names to match the directories on your
49 % computer. Furthermore, change the save file name to the name you want
50 % the save file to be.
51 save_file_str = [file_list{1}(1:end-4) '_vs_' file_list{2}(1:end-4)];
52 csv_dir = [pwd '/Data-Files/Dienstman_CSV-Files_' preprocessing '_PP'];
53 save_dir = [pwd '/Figure-Files/Dienstman_Univariate-Risks_' preprocessing '_PP'];
54 save_file = [save_dir '/' save_file_str '.png'];
55
56 % If the save_dir doesn't already exist, we make the save_dir here.
57 if ~exist(save_dir, 'dir')
58     mkdir(save_dir)
59 end
60
61 % If the save file already exists, we stop the program so we don't
62 % overwrite the file.
63 if exist(save_file, 'file')
64     disp('Error: A file already exists with the save file name. The program stopped because
65         running the program would overwrite the existing file.')
66     return
67 end
68
69 % These vectors are used to create the information in the boxes at
70 % the top of the figures.
71 N = length(file_list);
72 baby_count_vector = zeros(1, N);
73 mean_half_hour_vector = zeros(1, N);
74
75 % Here, we define the HRC names.
76 hrc_names = {'Asymmetry-1', 'Asymmetry-2', 'Asymmetry-Ratio', 'Decelerations', 'Mean-RR', '
77             'Sample-Entropy', 'Variance'};
78 hrc_types = {'', '_10', '_50', '_90', '_Slope'};
79 x_axis_list = {[0 10], [0 10], [0 5], [0 50], [250 600], [0 1.5], [0 5]};
80 x_axis_list_slope = {[ -5 5], [ -5 5], [ -1 1], [ -10 10], [ -50 50], [ -1 1], [ -2 2]};
81 x_axis_plot = {[ -2.5 10], [ -2.5 10], [ -1 4], [ -10 35], [ -60 560], [ -0.5 1.25], [ -1 4.5]};
82
83 % We take the header of one CSV file so we can look up the index of
84 % each HRC based on the string name.
85 [~, header_names] = xlsread([csv_dir '/' file_list{1}], '1:1');
86 site_index = find(strcmp(header_names, 'Site'));
87 id_index = find(strcmp(header_names, 'ID'));
88
89 % Next, we loop through every file in file_list and extract the
90 % necessary info.
91 for i = 1:N
92     csv_file = [csv_dir '/' file_list{i}];
93     raw_data = dlmread(csv_file, ',', 1, 0);
94     unique_matrix = unique(raw_data(:, [site_index id_index]), 'rows');
95     baby_count_vector(i) = length(unique_matrix);
96     mean_half_hour_vector(i) = length(raw_data)/length(unique_matrix);
97     eval(['raw_data_' num2str(i) ' = raw_data;'])
98 end
99
100 % Next, for each HRC, we calculate the risk. For a complete explanation
101 % on row to calculate the risk, see the documetation for
102 % one_risk_matrix.m.
103 for i = 1:length(hrc_names)
104     for j = 1:length(hrc_types)
105         variable_name = [hrc_names{i} hrc_types{j}];
106         variable_index = find(strcmp(header_names, variable_name)); %ok< *NASGU
107
108         if j ~= 5
109             xmin = x_axis_list{i}(1);
110             xmax = x_axis_list{i}(2);
111         else
112             xmin = x_axis_list_slope{i}(1);
113             xmax = x_axis_list_slope{i}(2);
114         end
115
116         if N == 2
117             eval([' variable_name '_Risk ' variable_name '_Points] = one_risk_matrix(
118                 preprocessing, variable_index, xmin, xmax, raw_data-1, raw_data-2;'])

```



```

117     else
118         eval([' variable_name '_Risk ' variable_name '_Points] = one_risk_matrix(
                preprocessing, variable_index, xmin, xmax, raw_data_1, raw_data_2, raw_data_3);'])
119     end
120 end
121 end
122
123 % Lastly, we graph the risks.
124 color_list = {[0 0 1], [0 1 0], [1 0 0], [0 1 1], [1 0 1]};
125 color_name_list = {'Blue', 'Green', 'Red', 'Cyan', 'Magenta'};
126 title_str = {'Raw', '10', '50', '90', 'Slope'};
127 probability_figure = figure('Position', [50,50,1600,900]);
128 set(probability_figure, 'color', 'w');
129
130 % Each subplot corresponds to one specific HRC type. On each subplot,
131 % there are five curves representing the risk for each type of HRC
132 % (raw, slope, 10, 50, and 90).
133 for i = 1:length(hrc_names)
134     subplot(3,3,i+1)
135     hold on
136
137     for j = 1:length(hrc_types)
138         variable_name = [hrc_names{i} hrc_types{j}];
139         eval(['x = ' variable_name '_Points;'])
140         eval(['y = ' variable_name '_Risk;'])
141         plot(x, y, 'LineWidth', 2, 'Color', color_list{j})
142     end
143
144     xlim(x_axis_plot{i})
145     ylim([-2 2])
146     plot(linspace(x_axis_plot{i}(1), x_axis_plot{i}(2), 100), zeros(1,100), '--', 'Color',
147         [0,0,0], 'LineWidth', 1);
148     xlabel(regexp(hrc_names{i}, '_', ''), 'FontSize', 16)
149     ylabel('Risk', 'FontSize', 16)
150     set(gca, 'FontSize', 16)
151     hold off
152 end
153
154 % Here, we plot the box at the top of the figure.
155 if strcmp(file_list{1}, 'hrc_sick.org.all_vent_all.csv')
156     baby_title_string = [ ' Univariate Risks: ' char(10)...
157         ' file_list{1}(1:end-4) ' vs. ' file_list{2}(1:end-4) char(10) ' '];
158
159     for n = 1:length(hrc_types)
160         baby_title_string = [baby_title_string color_name_list{n} ' = ' title_str{n} ' '; ']; %ok<*>
161         AGROW>
162     end
163
164     for n = 1:length(file_list)
165         file_str = file_list{n};
166         baby_title_string = [baby_title_string char(10) char(10)...
167             ' file_str(1:end-4) char(10)...
168             ' Baby Count = ' num2str(baby_count_vector(n)) ', Mean Half Hours = '
169             num2str(mean_half_hour_vector(n), 4)];
170     end
171 else
172     baby_title_string = [ ' Univariate Risks: ' char(10)...
173         ' file_list{1}(1:end-4) ' vs. ' file_list{2}(1:end-4) char(10) ' '];
174
175     for n = 1:length(hrc_types)
176         baby_title_string = [baby_title_string color_name_list{n} ' = ' title_str{n} ' '; '];
177     end
178
179     for n = 1:length(file_list)
180         file_str = file_list{n};
181         baby_title_string = [baby_title_string char(10) char(10)...
182             ' file_str(1:end-4) char(10)...
183             ' Baby Count = ' num2str(baby_count_vector(n)) ', Mean Half Hours = '
184             num2str(mean_half_hour_vector(n), 4)];
185     end
186 end
187
188 title_frame = uicontrol('style', 'frame');
189 set(title_frame, 'Position', [208, 635, 345, 195], 'BackgroundColor', [0 0 0])
190 baby_title = uicontrol('style', 'text');
191 set(baby_title, 'String', baby_title_string, 'Position', [211, 638, 339, 189], 'FontSize', 8, '
192     BackgroundColor', [1 1 1])
193
194 % Finally, we save the figure.
195 print(probability_figure, save_file, '-dpng')
196 end

```

Listing B.19: one_risk_matrix.m

```

1 function [risk, points, P_sigs_given_group_1, P_sigs_given_group_2] = one_risk_matrix(
    preprocessing, variables, xmin, xmax, half_hours_group_1, half_hours_group_2,
    half_hours_group_3)
2 % Author: Evan Dienstman
3 % Last Update: 3/31/2016
4 % Email: eddienstman@email.wm.edu
5 % Note: Feel free to email me with questions! If something doesn't

```

```

6 % make sense, it might be because I haven't updated the code yet.
7 %
8 % This function calculates the risk of being in different groups given
9 % one or two HRC signals. The function calculates the risk over many
10 % different values of the signals. The exact definition of risk is
11 % defined in the function. We calculate the risk using the data
12 % provided in the "half.hours.group" arguments. The specific HRC values
13 % in question are determined by the argument "variables". The function
14 % determines which HRC values used for the risk calculation.
15 % Thus, the user cannot query specific values to get their risks.
16 %
17 % Arguments:
18 % 1. preprocessing - a string indicating which preprocessing method
19 % to use
20 % 2. variables - a vector of number(s) that indicate which HRC
21 % signal(s) to use when calculating the risk
22 % 3. xmin - the minimum HRC value to evaluate the risk at
23 % 4. xmax - the maximum HRC value to evaluate the risk at
24 % 5. half.hours.group_1 - the HRC data for group 1 (most likely the
25 % sick half hours or sick organism x half hours)
26 % 6. half.hours.group_2 - the HRC data for group 2 (most likely the
27 % healthy half hour group)
28 % 7. half.hours.group_3 - the HRC data for group 3 (most likely the
29 % sick organism y group)
30 %
31 % Preconditions:
32 % 1. Make sure the files one_porb_matrix.m hrc_indices.m, and
33 % bin_width_ma.m are in the working directory.
34 %
35 % Returns:
36 % 1. risk - the risk at each value in the meshgrid (or vector for
37 % for one signal) of the variable "points"
38 % 2. points - the matrix (or vector for one signal) of various
39 % HRC values for the signal(s) specified where the meshgrid of
40 % the rows in "points" correspond to the location of each "risk"
41 % value
42 % 3. P_sigs_given_group_1 - the probability at each value in the
43 % meshgrid of the variable "points" for the signal(s) given
44 % group 1
45 % 4. P_sigs_given_group_2 - the probability at each value in the
46 % meshgrid of the variable "points" for the signal(s) given
47 % group 2
48 %
49 % First, we calculate the number of half hours we have.
50 num_half_hours_group_1 = length(half_hours_group_1);
51 num_half_hours_group_2 = length(half_hours_group_2);
52 total_half_hours = num_half_hours_group_1 + num_half_hours_group_2;
53 %
54 % Next, we find the probability of being in group 1, probability of
55 % being in group 2, probability of getting various signals given
56 % group 1, probability of getting various signals given group 2, and
57 % probability of getting various signals. Notice that P_group_1 and
58 % P_group_2 are numbers while the other probabilities are matrices.
59 % The mesh grid of the rows in "points" indicate the signal values(s)
60 % of the probability matrices.
61 P_group_1 = num_half_hours_group_1/total_half_hours;
62 P_group_2 = num_half_hours_group_2/total_half_hours;
63 [P_sigs_given_group_1, points] = one_prob_matrix(variables, half_hours_group_1, xmin, xmax,
64 preprocessing);
65 P_sigs = P_sigs_given_group_1 .* P_group_1 + P_sigs_given_group_2 .* P_group_2;
66 %
67 % We define the risk in two different ways depending on if we are
68 % comparing sick to healthy half hours or sick organism x to sick
69 % organism y half hours.
70 if ~exist('half_hours_group_3', 'var')
71 %
72 % Here, we define the risk for comparing sick half hours to
73 % healthy half hours. Note that the sick half hours can be either
74 % group 1 or group 2, but generally we will use group 1 for the
75 % sick half hours. Therefore, values above 0 indicates the patient
76 % is more likely to be sick. We also remove any risks at points
77 % where the probabilities used to calculate the risk is very small
78 % (corresponding to little data).
79 risk = log(P_sigs_given_group_1 ./ P_sigs);
80 risk(P_sigs_given_group_1 < 0.001 & P_sigs_given_group_2 < 0.001) = NaN;
81 %
82 else
83 %
84 % Here, we define the risk for comparing sick organism x to sick
85 % organism y. In short, we first calculate the risk of sick
86 % organism x to healthy and sick organism y to healthy. Afterwards,
87 % we take the ratio of these two risks to get a new risk. Thus, we
88 % can interpret values above 0 as more likely to be sick from
89 % organism x and values below zero as more likely to be sick from
90 % organism y. Again, we remove risks at points where the
91 % probabilities are very small.
92 num_half_hours_group_3 = length(half_hours_group_3);
93 total_half_hours_new = num_half_hours_group_3 + num_half_hours_group_2;
94 P_group_3 = num_half_hours_group_3/total_half_hours_new;
95 P_group_2_new = num_half_hours_group_2/total_half_hours_new;
96 P_sigs_given_group_3 = one_prob_matrix(variables, half_hours_group_3, xmin, xmax,
97 preprocessing);
98 P_sigs_given_group_2_new = one_prob_matrix(variables, half_hours_group_2, xmin, xmax,

```

```

    preprocessing);
98     P_sigs_new = P_sigs_given_group_3 * P_group_3 + P_sigs_given_group_2_new * P_group_2_new;
99
100    risk1 = P_sigs_given_group_1./P_sigs;
101    risk3 = P_sigs_given_group_3./P_sigs_new;
102    risk = log(risk1./risk3);
103    risk( (P_sigs_given_group_1 < 0.001 & P_sigs_given_group_2 < 0.001) | (P_sigs_given_group_3 <
        0.001 & P_sigs_given_group_2_new < 0.001) ) = NaN;
104 end
105 end

```

Listing B.20: one_prob_matrix.m

```

1  function [prob, points] = one_prob_matrix(variables, data_matrix, xmin, xmax, preprocessing)
2  % Author: Evan Dienstman
3  % Last Update: 3/31/2016
4  % Email: eddienstman@email.wm.edu
5  % Note: Feel free to email me with questions!
6  %
7  % This function calculates the probability of getting various HRC
8  % values using the data provided in "data_matrix". The function
9  % determines which HRC values used for the probability calculation.
10 % Thus, the user cannot query specific values to get their probability.
11 % The specific HRCs in question are determined by the argument
12 % "variables". The argument "variables" can only contain one or two
13 % numbers corresponding to the probability of HRC X or the bivariate
14 % probability of HRC X and HRC Y. To calculate the probability, we use
15 % kernel density estimation to get a PDF, and then we integrate the PDF
16 % over a binwidth to get the probability.
17 %
18 % Arguments:
19 % 1. variables - a vector of number(s) that indicate which HRC
20 %    signal(s) to use from data_matrix
21 % 2. data_matrix - the data used to calculate the probability of
22 %    various values of the HRC signal(s) specified
23 % 3. xmin - the minimum HRC value to evaluate the probability at
24 % 4. xmax - the maximum HRC value to evaluate the probability at
25 % 5. preprocessing - a string indicating which preprocessing method
26 %    to use
27 %
28 % Preconditions:
29 % 1. The Matlab function ksdensity for bivariate PDFs only works
30 %    on Matlab 2016 or later.
31 % 2. Make sure the file hrc_indices.m and bin_widths.m are in the
32 %    working directory.
33 %
34 % Returns:
35 % 1. prob - the probability at each value in variable "points"
36 % 2. points - the matrix (or vector from one signal) of various HRC
37 %    values for the signal(s) specified where the meshgrid of the
38 %    rows in "points" correspond to the location of each "prob"
39 %    value
40 %
41 % First, we load the binwidths we use for the bandwidth of ksdensity and
42 % when integrating the PDF to calculate the probability. For more
43 % information of how we calculate the binwidth, please see the
44 % documentation for csv_bin_widths.m.
45 load(['bin_widths_' preprocessing '_PP.mat']);
46 load('hrc_indices.mat')
47
48 % We must split up the function into two cases: 1) when we want the
49 % probability given two signals and 2) when we want the probability
50 % given one signal.
51 if length(variables) == 2
52
53     % Here, we find the probability for a matrix of (x,y)-points given
54     % two HRC signals. We use N = 32 because a 32x32 matrix gives us a
55     % good plot for not too much computation time.
56     N = 32;
57
58     % Next, we prepare the data for the ksdensity function.
59     var1 = variables(1);
60     var2 = variables(2);
61
62     x_points = linspace(xmin(1), xmax(1), N);
63     y_points = linspace(xmin(2), xmax(2), N);
64     points = [x_points; y_points];
65     [X, Y] = meshgrid(x_points, y_points);
66
67     data1 = data_matrix(:, var1);
68     data2 = data_matrix(:, var2);
69
70     prob = zeros(N,N);
71
72     % Here, we visit each point the the matrix one by one and the
73     % calculate the probability around that point.
74     for i = 1:N
75
76         for j = 1:N
77             x = X(i, j);
78             y = Y(i, j);
79             local_x_points = [x-0.5*bin_widths(var1) x+0.5*bin_widths(var1)];

```

```

80     local_y_points = [y-0.5*bin_widths(var2) y+0.5*bin_widths(var2)];
81     [local_X, local_Y] = meshgrid(local_x_points, local_y_points);
82
83     % Next, we calculate the bivariate PDF using the Matlab
84     % function ksdensity. We do not let ksdensity pick an
85     % optimal bandwidth because the bandwidth for
86     % decelerations should be 1. Since there is no way to set
87     % the bandwidth for decelerations as 1 and generate the
88     % other bandwidth automatically, we simply use our binwidths
89     % as the bandwidth. We also use the Epanechnikov kernel.
90     % For more information on kernel density estimation, please
91     % refer to its Wikipedia page.
92     p = ksdensity([data1 data2], [local_X(:) local_Y(:)], 'kernel', 'epanechnikov', 'width
93     ', [0.5*bin_widths(var1) 0.5*bin_widths(var2)]);
94     p = reshape(p, size(local_X));
95
96     % Finally, we integrate the PDF to get the probability in a
97     % certain area. Notice that we use our binwidth for the
98     % local_x_points and local_y_points. Thus, when we
99     % integrate, we expect a value close to the true
100    % probability in that area.
101    prob(i,j) = trapz( local_y_points, trapz(local_x_points,p,2) );
102
103    end
104
105    else
106
107    % Here, we find the probability for a vector of x-points given one
108    % HRC signal. The x-point vector is defined below. Note that each
109    % x-point is one "width" apart. The importance of this feature
110    % is explained below.
111    data = data_matrix(:, variables);
112    x_points = (xmin:bin_widths(variables):xmax).';
113    points = x_points(2:end) - 0.5*bin_widths(variables);
114    N = (length(x_points)-1);
115    prob = zeros(N,1);
116
117    % Next, we calculate the PDF using the Matlab function ksdensity.
118    % Here, we let ksdensity chose the optimal bandwidth except for the
119    % case of decelerations where the width should be 1. We also use
120    % the Epanechnikov kernel. For more information on kernel
121    % density estimation, please refer to its Wikipedia page.
122    if variables >= decelerations_index && variables <= decelerations_90_index
123        p = ksdensity(data, x_points, 'kernel', 'epanechnikov', 'width', 1);
124    else
125        p = ksdensity(data, x_points, 'kernel', 'epanechnikov');
126    end
127
128    % Next, we integrate the PDF over a standard bandwidth to get the
129    % probabilities around certain points. We use a slightly
130    % different width than what ksdensity chose above so all the
131    % vectors for like HRC's contain the same number of points.
132    % In this manner, we can add these vectors later in
133    % one_risk_matrix(). Also, since we step by a "width" in the
134    % x-points vector, integrating between consecutive points will give
135    % us a probability around the points close to the true value.
136    for i = 1:N
137        prob(i) = trapz([x_points(i) x_points(i+1)], [p(i) p(i+1)]);
138    end
139    end
140
141    end
142
143    end

```

B.6 Bivariate PDF Figures

Listing B.21: multiple_bivariate_pdf_figures.m

```

1 % Author: Evan Dienstman
2 % Last Update: 3/31/2016
3 % Email: eddienstman@email.wm.edu
4 % Note: Feel free to email me with questions! If something doesn't
5 % make sense, it might be because I haven't updated the code yet.
6 %
7 % This script calls the function one_bivariate_pdf_maker using
8 % different combinations of csv files. Each call to the function is
9 % responsible for one figure found in the folder
10 % Dienstman_Bivariate_PDFs. For more information on what is in these
11 % figures and how they were created, see the documentations for
12 % one_bivariate_pdf_maker.
13 %
14 % Preconditions:
15 %     1. Make sure the files one_bivariate_pdf_figure.m, hrc_indces,
16 %     and bin_width.m are in the working directory.
17 %     2. Make sure to select the proper preprocessing method below.
18 %
19 % Returns:
20 %     1. This function creates one figure for every call to
21 %     one_bivariate_pdf_figure saved in the folder

```

```

22 %      Dienstman_bivariate_PDFs.
23
24 clear
25 clc
26
27 % Change the variable preprocessing to match the preprocessing method
28 % used to make the CSV files that you will use below to create the
29 % bivariate PDFs.
30 preprocessing = 'Abby';
31 % preprocessing = 'Doug';
32
33 % These are the column numbers of the CSV files we will use when
34 % making different bivariate PDFs. Each column corresponds to a
35 % different HRC. Change these numbers if you wish to use other HRCs.
36 % Note that we only use the raw HRCs because including all subcategories
37 % would create too many bivariate PDFs. We also don't use mean RR,
38 % asymmetry 1, and asymmetry 2 for the same reason.
39 load('hrc_indices.mat')
40 variables = [variance_index sample_entropy_index asymmetry_ratio_index decelerations_index];
41
42 % Here, we call one_bivariate_pdf_maker with different combinations
43 % of CSV files. The function one_bivariate_pdf_maker will create the
44 % figures saved in the folder Dienstman_Bivariate_PDFs using the two
45 % CSV files you pass into the function.
46 one_bivariate_pdf_figure(variables, 'hrc_healthy_org_all_vent_all.csv', 'hrc_sick_org_all_vent_all
47 .csv', preprocessing)
48 one_bivariate_pdf_figure(variables, 'hrc_healthy_org_all_vent_1.csv', 'hrc_sick_org_all_vent_1.csv
49 ', preprocessing)
50 one_bivariate_pdf_figure(variables, 'hrc_healthy_org_all_vent_0.csv', 'hrc_sick_org_all_vent_0.csv
51 ', preprocessing)
52 one_bivariate_pdf_figure(variables, 'hrc_sick_org_1_vent_0.csv', 'hrc_sick_org_2_vent_0.csv',
53 preprocessing)
54 one_bivariate_pdf_figure(variables, 'hrc_sick_org_1_vent_0.csv', 'hrc_sick_org_3_vent_0.csv',
55 preprocessing)
56 one_bivariate_pdf_figure(variables, 'hrc_sick_org_1_vent_0.csv', 'hrc_sick_org_4_vent_0.csv',
57 preprocessing)
58 one_bivariate_pdf_figure(variables, 'hrc_sick_org_1_vent_0.csv', 'hrc_sick_org_5_vent_0.csv',
59 preprocessing)
60 one_bivariate_pdf_figure(variables, 'hrc_sick_org_2_vent_0.csv', 'hrc_sick_org_3_vent_0.csv',
61 preprocessing)
62 one_bivariate_pdf_figure(variables, 'hrc_sick_org_2_vent_0.csv', 'hrc_sick_org_4_vent_0.csv',
63 preprocessing)
64 one_bivariate_pdf_figure(variables, 'hrc_sick_org_2_vent_0.csv', 'hrc_sick_org_5_vent_0.csv',
65 preprocessing)
66 one_bivariate_pdf_figure(variables, 'hrc_sick_org_3_vent_0.csv', 'hrc_sick_org_4_vent_0.csv',
67 preprocessing)
68 one_bivariate_pdf_figure(variables, 'hrc_sick_org_3_vent_0.csv', 'hrc_sick_org_5_vent_0.csv',
69 preprocessing)
70 one_bivariate_pdf_figure(variables, 'hrc_sick_org_4_vent_0.csv', 'hrc_sick_org_5_vent_0.csv',
71 preprocessing)

```

Listing B.22: one_bivariate_pdf_figure.m

```

1 function one_bivariate_pdf_figure(variables, file1, file2, preprocessing)
2 % Author: Evan Dienstman
3 % Last Update: 3/31/2016
4 % Email: eddiendienstman@email.wm.edu
5 % Note: Feel free to email me with questions! If something doesn't
6 % make sense, it might be because I haven't updated the code yet.
7 %
8 % This function creates bivariate probability density functions (PDFs)
9 % for certain combinations of two heart rate characteristics (HRCs).
10 % The PDFs are then all plotted onto one figure. There are six subplots
11 % in the figure, one for each combination of HRCs in the input argument
12 % "variables", and each subplot has two PDFs that corresponds to file1
13 % and file2. Note that the argument "variables" must contain exactly
14 % four indices (making six combinations) because that the figure is
15 % only set up to make six subplots. There is also a subtitle at the top
16 % of the figure that acts as a legend for each subplot. Note that for
17 % any sick category, we take a period of 24 half hours from an event
18 % and for any healthy category, we take a period of 7 days before and
19 % event or 3 days after an event. The PDFs are created using kernel
20 % density estimation. For more details about this method, see the
21 % comments in the code.
22 %
23 % Arguments:
24 % 1. variables - the six column numbers from the CSV files to use
25 % when making the bivariate PDFs
26 % 2. file1 - the first file to calculate the bivariate PDFs of for
27 % each combination of the indices in variables
28 % 3. file2 - the second file to calculate the bivariate PDFs of for
29 % each combination of the indices in variables
30 % 4. preprocessing - the preprocessing method used when determining
31 % which CSV files to use
32 %
33 % Preconditions:
34 % 1. The Matlab function ksdensity for bivariate PDFs only works
35 % on Matlab 2016 or later.
36 % 2. Make sure the directories and file names used in the scripts
37 % are the right ones for the computer you are using.
38 % 3. This script will not overwrite any existing files with the
39 % same name. Change the variable save_file to a name that does

```

```

40 %         not already exist, or delete the existing file with the same
41 %         name before running this script.
42 %         4. Make sure the file bin_widths.m is in the working
43 %         directory.
44 %         5. The input argument "variables" must contain exactyl four
45 %         indices.
46 %
47 % Returns:
48 %     1. This function returns a figure saved with the name given by
49 %        the variable save_file that shows the bivariate PDFs for each
50 %        combination of HRCs in variables for file1 and file2.
51 %
52 % Change the directory names to match the directories on your
53 % computer. Furthermore, change the bin width file name to the name of
54 % the file already on your computer and change the save file name to
55 % the name you want the save file to be.
56 csv_dir = [pwd '/Data-Files/Dienstman_CSV-Files_' preprocessing '_PP'];
57 save_dir = [pwd '/Figure-Files/Dienstman-Bivariate-PDFs_' preprocessing '_PP'];
58 save_file = [save_dir '/' file1(1:end-4) '_vs_' file2(1:end-4) '.png'];
59 bin_width.file = ['bin_widths_' preprocessing '_PP.mat'];
60
61 % If the save_dir doesn't already exist, we make the save_dir here.
62 if ~exist(save_dir, 'dir')
63     mkdir(save_dir)
64 end
65
66 % If the save file already exists, we stop the program so we don't
67 % overwrite the file.
68 if exist(save_file, 'file')
69     disp('Error: A file already exists with the save file name. The program stopped because
70         running the program would overwrite the existing file.')
71     return
72 end
73
74 % Here, we prepare the title for our figure by using the names from
75 % the csv files given.
76 title_name_1 = file1(1:end-4);
77 title_name_1 = regexp(title_name_1, '_', ' ');
78
79 title_name_2 = file2(1:end-4);
80 title_name_2 = regexp(title_name_2, '_', ' ');
81
82 % Here, we open up the CSV files. Change the file name to match the
83 % location of the file on your computer.
84 raw_data_1 = dlmread([csv_dir '/' file1], ',', 1, 0);
85 raw_data_2 = dlmread([csv_dir '/' file2], ',', 1, 0);
86
87 % Next, we load in the bivariate bin widths used for smoothing and
88 % creating the edges for the PDFs.
89 load(bin_width_file)
90
91 % Here, we load the header. We use the smallest file since all the
92 % headers are the same.
93 [~, variable_names] = xlsread([csv_dir '/' hrc_sick_org_5_vent_0.csv'], '1:1');
94
95 % We use N = 32 because a 32x32 matrix gives us a good plot for not too
96 % much computation time.
97 N = 32;
98
99 % Next, we loop through each plot for each file (six plots for two
100 % files for a total of 12 plots) and create the bivariate PDFs for each
101 % plot.
102 for i = 1:12
103     % If i <= 6, we are plotting the first file's PDFs. Else, we are
104     % plotting the second file's PDFs.
105     if i <= 6
106         raw_data = raw_data_1;
107         plot_num = i;
108     else
109         raw_data = raw_data_2;
110         plot_num = i - 6;
111     end
112
113     % If i == 1, we create the figure and the title for the figure.
114     if i == 1
115         hist_fig = figure('Position', [50,50,1600,900]);
116         set(hist_fig, 'color', 'w');
117         title_string = ['Bivariate Probability Densities: ' title_name_1 ' and ' title_name_2];
118         sub_title_string = ['blue = ' title_name_1 ' ; red = ' title_name_2];
119         figure_title = uicontrol('style', 'text');
120         figure_sub_title = uicontrol('style', 'text');
121         set(figure_title, 'String', title_string, 'Position', [325, 860, 1000, 40], 'FontSize',
122             18, 'BackgroundColor', [1 1 1])
123         set(figure_sub_title, 'String', sub_title_string, 'Position', [590, 830, 500, 30], '
124             FontSize', 12, 'BackgroundColor', [1 1 1])
125     end
126
127 % Below, we check which iteration we are on to determine which
128 % subplot we are at. For each subplot, we must determine the HRC
129 % for the x_edge and y_dge of the bivariate PDF. Looking at the
130 % code below, the variables x_edge and y_edge can be any number in
131 % "variables" corresponding to an HRC column in the data file.
132 if plot_num <= 3

```

```

131     x_edge = variables(1);
132     x_plot_edge = linspace(0,5,N);
133     y_edge = variables(plot_num + 1);
134
135     % Change the axes for the HRCs you are using.
136     if y_edge == variables(2)
137         y_plot_edge = linspace(0,1,N);
138         z_axis = [0 3];
139     elseif y_edge == variables(3)
140         y_plot_edge = linspace(0,2.5,N);
141         z_axis = [0 2];
142     elseif y_edge == variables(4)
143         y_plot_edge = linspace(0,10,N);
144         z_axis = [0 0.5];
145     end
146
147     elseif plot_num <= 5
148         x_edge = variables(2);
149         x_plot_edge = linspace(0,1,N);
150         y_edge = variables(plot_num - 1);
151
152         % Change the axes for the HRCs you are using.
153         if y_edge == variables(3)
154             y_plot_edge = linspace(0,2.5,N);
155             z_axis = [0 10];
156         elseif y_edge == variables(4)
157             y_plot_edge = linspace(0,10,N);
158             z_axis = [0 1.5];
159         end
160
161     elseif plot_num == 6
162         x_edge = variables(3);
163         x_plot_edge = linspace(0,2.5,N);
164         y_edge = variables(4);
165
166         % Change the axes for the HRCs you are using.
167         y_plot_edge = linspace(0,10,N);
168         z_axis = [0 1];
169     end
170
171     % The variables x_data and y_data are the two data vecotrs used to
172     % make the PDFs. Once we have selected the two data vectors, we can
173     % create the bivariate PDFs. We create PDFs using the Matlab
174     % function ksdensity, which uses the kernel density estimation
175     % method to create the PDFs. For the kernel density estimation, we
176     % use an Epanechnikov kernel and band widths equal to the bin widths
177     % of each HRC. We use the bin widths we calculated earlir so all
178     % the PDFs are consistent with the band widths used. We feel we can
179     % more justly compare PDFs in this manner. For more information
180     % about our choice of bin widths, see the documentation for
181     % bin-width-maker. We also evaluate the PDFs at the edges
182     % calculated earlier. The input for ksdensity is a two column
183     % matrix with each column corresponding to one specfic type of HRC
184     % values. Each row in the matrix corresponds to every half hour in
185     % the csv file in question. Note that since we are creating PDFs,
186     % the volume underneath each surface is 1. Consequently, all the
187     % PDFs are normalized and values of the kernel PDF can be greater than 1.
188     % For more information about the kernel density estimation,
189     % Epanechnikov kernel function, and PDFs, see their respective
190     % wikipedia pages.
191     x_data = raw_data(:,x_edge);
192     y_data = raw_data(:,y_edge);
193     [X, Y] = meshgrid(x_plot_edge, y_plot_edge);
194     Z = ksdensity([x_data y_data], [X(:) Y(:)], 'kernel', 'epanechnikov', 'width', [0.5*bin_widths
195         (x_edge) 0.5*bin_widths(y_edge)]);
196     Z = reshape(Z, size(X));
197
198     % Lastly, we can plot the bivariate PDFs. We plot the first
199     % file's PDFs in blue and the second file's PDFs in red.
200     hold on
201     subplot(3,2,plot_num)
202
203     if i <= 6
204         surf(X, Y, Z, 'facealpha', .5, 'Facecolor', 'b')
205     else
206         surf(X, Y, Z, 'facealpha', .5, 'Facecolor', 'r')
207     end
208
209     xlabel(regexprep(variable_names{x_edge}, '_',' '), 'FontSize', 14)
210     ylabel(regexprep(variable_names{y_edge}, '_',' '), 'FontSize', 14)
211     zlabel('Probability', 'FontSize', 14)
212     set(gca, 'FontSize', 14)
213     xlim([x_plot_edge(1) x_plot_edge(end)])
214     ylim([y_plot_edge(1) y_plot_edge(end)])
215     zlim(z_axis)
216     caxis(z_axis)
217     hold off
218 end
219
220 % Finally, we save the figure.
221 print(hist_fig, save_file, '-dpng')
222 end

```

B.7 Bivariate Risk Figures

Listing B.23: multiple_bivariate_risk_figures.m

```
1 % Author: Evan Dienstman
2 % Last Update: 3/31/2016
3 % Email: eddiendienstman@email.wm.edu
4 % Note: Feel free to email me with questions! If something doesn't
5 % make sense, it might be because I haven't updated the code yet.
6 %
7 % This script creates multiple bivariate risk figures. Each figure
8 % calculates the risk associated with different half hour groups for
9 % many different combination of two HRCs. For more information about
10 % the figures, individual plots, and how we define the risk, please see
11 % the documentation for one_bivariate_risk_figure.m,
12 % one_bivariate_risk_plot.m, and one_risk_matrix.m.
13 %
14 % Preconditions:
15 % 1. Make sure the files one_bivariate_risk_figure.m,
16 % one_bivariate_risk_plot.m, one_risk_matrix.m,
17 % one_porb_matrix.m, bin_widths.m, and hrc_indices.m are in the
18 % working directory.
19 % 2. Make sure to select the proper preprocessing method below.
20 %
21 % Returns:
22 % 1. This function returns one bivariate risk figure for every call
23 % to one_bivariate_risk_figure.m. For info on where the figures
24 % are saved, see the documentation of one_bivariate_risk.m.
25 %
26 clear
27 clc
28
29 % Here, we select the preprocessing method.
30 preprocessing = 'Abby';
31 preprocessing = 'Doug';
32
33 % These are the column numbers of the CSV files we will use when
34 % making different bivariate risks. Each column corresponds to a
35 % different HRC. Change these numbers if you wish to use other HRCs.
36 % Note that we only use the raw HRCs because including all subcategories
37 % would create too many bivariate risks. We also don't use mean
38 % RR, asymmetry 1, and asymmetry 2 for the same reason.
39 load('hrc_indices.mat')
40 variables = [variance_index sample_entropy_index asymmetry_ratio_index decelerations_index];
41
42 % Here we define the xmin and xmax for plotting.
43 xmin_vector = {[1 0], [1 0.5], [1 0], [0 0.5], [0 0], [0.5 0]};
44 xmax_vector = {[4 0.8], [4 2.5], [4 10], [0.8 2.5], [0.8 10], [2.5 10]};
45
46 % Next, we call the procedure to create each figure. Each figure
47 % uses a different combination of half hour groups.
48 % one_bivariate_risk_figure(preprocessing, variables, xmin_vector, xmax_vector, '
49 % hrc_sick_org_all_vent_all.csv', 'hrc_healthy_org_all_vent_all.csv')
50 % one_bivariate_risk_figure(preprocessing, variables, xmin_vector, xmax_vector, '
51 % hrc_sick_org_all_vent_1.csv', 'hrc_healthy_org_all_vent_1.csv')
52 % one_bivariate_risk_figure(preprocessing, variables, xmin_vector, xmax_vector, '
53 % hrc_sick_org_all_vent_0.csv', 'hrc_healthy_org_all_vent_0.csv')
54 % one_bivariate_risk_figure(preprocessing, variables, xmin_vector, xmax_vector, '
55 % hrc_sick_org_1_vent_0.csv', 'hrc_sick_org_all_vent_0.csv', 'hrc_sick_org_2_vent_0.csv')
56 % one_bivariate_risk_figure(preprocessing, variables, xmin_vector, xmax_vector, '
57 % hrc_sick_org_1_vent_0.csv', 'hrc_sick_org_all_vent_0.csv', 'hrc_sick_org_3_vent_0.csv')
58 % one_bivariate_risk_figure(preprocessing, variables, xmin_vector, xmax_vector, '
59 % hrc_sick_org_1_vent_0.csv', 'hrc_sick_org_all_vent_0.csv', 'hrc_sick_org_4_vent_0.csv')
60 % one_bivariate_risk_figure(preprocessing, variables, xmin_vector, xmax_vector, '
61 % hrc_sick_org_2_vent_0.csv', 'hrc_sick_org_all_vent_0.csv', 'hrc_sick_org_5_vent_0.csv')
62 % one_bivariate_risk_figure(preprocessing, variables, xmin_vector, xmax_vector, '
63 % hrc_sick_org_2_vent_0.csv', 'hrc_sick_org_all_vent_0.csv', 'hrc_sick_org_3_vent_0.csv')
64 % one_bivariate_risk_figure(preprocessing, variables, xmin_vector, xmax_vector, '
65 % hrc_sick_org_2_vent_0.csv', 'hrc_sick_org_all_vent_0.csv', 'hrc_sick_org_4_vent_0.csv')
66 % one_bivariate_risk_figure(preprocessing, variables, xmin_vector, xmax_vector, '
67 % hrc_sick_org_2_vent_0.csv', 'hrc_sick_org_all_vent_0.csv', 'hrc_sick_org_5_vent_0.csv')
68 % one_bivariate_risk_figure(preprocessing, variables, xmin_vector, xmax_vector, '
69 % hrc_sick_org_3_vent_0.csv', 'hrc_sick_org_all_vent_0.csv', 'hrc_sick_org_4_vent_0.csv')
70 % one_bivariate_risk_figure(preprocessing, variables, xmin_vector, xmax_vector, '
71 % hrc_sick_org_3_vent_0.csv', 'hrc_sick_org_all_vent_0.csv', 'hrc_sick_org_5_vent_0.csv')
72 % one_bivariate_risk_figure(preprocessing, variables, xmin_vector, xmax_vector, '
73 % hrc_sick_org_4_vent_0.csv', 'hrc_sick_org_all_vent_0.csv', 'hrc_sick_org_5_vent_0.csv')
```

Listing B.24: one_bivariate_risk_figure.m

```
1 function one_bivariate_risk_figure(preprocessing, variables, xmin_vector, xmax_vector, file1,
2 file2, file3)
3 % Author: Evan Dienstman
4 % Last Update: 4/5/2016
5 % Email: eddiendienstman@email.wm.edu
6 % Note: Feel free to email me with questions! If something doesn't
7 % make sense, it might be because I haven't updated the code yet.
8 %
9 % This function creates a figure containing multiple risk plots.
```



```

9 % Each figure contains six subplots corresponding to six different
10 % combinations of the four input HRCs. Note that the input argument
11 % "variables" must contain exactly four HRCs because the figure is only
12 % set up for four. Each subplot then contains the appropriate risk plot
13 % for those two HRCs. The subplots are created by
14 % one_bivariate_risk_plot.m. For more information about the subplots
15 % and how we define the risk, please see the documentation for
16 % one_bivariate_risk_plot.m.
17 %
18 % Arguments:
19 % 1. preprocessing - a string indicating which preprocessing method
20 % to use
21 % 2. variables - the column indices from the CSV files of the four
22 % HRCs we want to use for the bivariate risks
23 % 3. xmin_vector - the minimum x-values for the subplots
24 % 4. xmax_vector - the maximum x-values for the subplots
25 % 5. file1 - the file containing the half hours for group 1
26 % 6. file2 - the file containing the half hours for group 2
27 % 7. file3 - the file containing the half hours for group 3
28 %
29 % Preconditions:
30 % 1. Make sure the files one_bivariate_risk_plot.m,
31 % one_risk_matrix.m, one_porb_matrix.m, and bin_width_maker.m
32 % are in the working directory.
33 % 2. Make sure there are exactly four indices in the input argument
34 % "variables".
35 % 3. Make sure xmin_vector and xmax_vector are of the form
36 % {[x1 y1], ..., [x4 y4]}.
37 % 4. This script will not overwrite any existing files with the
38 % same name. Change the variable save_file_str to a name that
39 % does not already exist, or delete the existing file with the
40 % same name before running this script.
41 %
42 % Returns:
43 % 1. This function returns a figure containing the six subplots
44 % produced by one_bivariate_risk_plot.m. Each plot represents
45 % the risk associate with a different combination of two HRCs
46 % relative to the two groups in question.
47 %
48 % First, we create the save directory.
49 save_dir = [pwd '/Figure-Files/Dienstman-Bivariate-Risks_' preprocessing '_PP'];
50
51 % If the save_dir doesn't already exist, we make the save_dir here.
52 if ~exist(save_dir, 'dir')
53     mkdir(save_dir)
54 end
55
56 % Next, we create the save file.
57 if exist('file3', 'var')
58     save_file_str_short = [file1(1:end-4) '_vs_' file3(1:end-4)];
59     save_file_str = [save_dir '/' file1(1:end-4) '_vs_' file3(1:end-4) '.png'];
60 else
61     save_file_str_short = [file1(1:end-4) '_vs_' file2(1:end-4)];
62     save_file_str = [save_dir '/' file1(1:end-4) '_vs_' file2(1:end-4) '.png'];
63 end
64
65 % If the save file already exists, we stop the program so we don't
66 % overwrite the file.
67 if exist(save_file_str, 'file')
68     disp('Error: A file already exists with the save file name. The program stopped because
69         running the program would overwrite the existing file.')
70     return
71 end
72
73 % Next, we do all the formatting for the figure.
74 title_str = ['Relative risk plots: ' regexprep(save_file_str_short, '_', ' ')];
75 sub_title_str = ['Values above 0 mean more likely to be ' regexprep(file1(1:end-4), '_', ' ') '.'];
76
77 plot_fig = figure('Position', [50,50,1600,900]);
78 set(plot_fig, 'color', 'w');
79 figure_title = uicontrol('style', 'text');
80 figure_sub_title = uicontrol('style', 'text');
81 set(figure_title, 'String', title_str, 'Position', [410, 845, 800, 40], 'FontSize', 18, '
82     BackgroundColor', [1 1 1])
83 set(figure_sub_title, 'String', sub_title_str, 'Position', [540, 830, 575, 20], 'FontSize', 12, '
84     BackgroundColor', [1 1 1])
85
86 % Here, we call one_bivariate_risk_plot.m with each combination of
87 % the HRCs to make the six subplots.
88 if exist('file3', 'var')
89     subplot(3,2,1)
90     one_bivariate_risk_plot(preprocessing, variables(1), variables(2), xmin_vector{1}, xmax_vector
91         {1}, file1, file2, file3)
92     subplot(3,2,2)
93     one_bivariate_risk_plot(preprocessing, variables(1), variables(3), xmin_vector{2}, xmax_vector
94         {2}, file1, file2, file3)
95     subplot(3,2,3)
96     one_bivariate_risk_plot(preprocessing, variables(1), variables(4), xmin_vector{3}, xmax_vector
97         {3}, file1, file2, file3)
98     subplot(3,2,4)
99     one_bivariate_risk_plot(preprocessing, variables(2), variables(3), xmin_vector{4}, xmax_vector
100        {4}, file1, file2, file3)
101     subplot(3,2,5)
102     one_bivariate_risk_plot(preprocessing, variables(2), variables(4), xmin_vector{5}, xmax_vector

```

```

    {5}, file1, file2, file3)
95 subplot(3,2,6)
96 one_bivariate_risk_plot(preprocessing, variables(3), variables(4), xmin_vector{6}, xmax_vector
    {6}, file1, file2, file3)
97
98 else
99 subplot(3,2,1)
100 one_bivariate_risk_plot(preprocessing, variables(1), variables(2), xmin_vector{1}, xmax_vector
    {1}, file1, file2)
101 subplot(3,2,2)
102 one_bivariate_risk_plot(preprocessing, variables(1), variables(3), xmin_vector{2}, xmax_vector
    {2}, file1, file2)
103 subplot(3,2,3)
104 one_bivariate_risk_plot(preprocessing, variables(1), variables(4), xmin_vector{3}, xmax_vector
    {3}, file1, file2)
105 subplot(3,2,4)
106 one_bivariate_risk_plot(preprocessing, variables(2), variables(3), xmin_vector{4}, xmax_vector
    {4}, file1, file2)
107 subplot(3,2,5)
108 one_bivariate_risk_plot(preprocessing, variables(2), variables(4), xmin_vector{5}, xmax_vector
    {5}, file1, file2)
109 subplot(3,2,6)
110 one_bivariate_risk_plot(preprocessing, variables(3), variables(4), xmin_vector{6}, xmax_vector
    {6}, file1, file2)
111 end
112
113 % Finally, we save the figure.
114 print(plot_fig, save_file_str, '-dpng')
115 end

```

Listing B.25: one_bivariate_risk_plot.m

```

1 function one_bivariate_risk_plot( preprocessing, var1, var2, xmin, xmax, file1, file2, file3)
2 % Author: Evan Dienstman
3 % Last Update: 3/31/2016
4 % Email: eddienstman@email.wm.edu
5 % Note: Feel free to email me with questions! If something doesn't
6 % make sense, it might be because I haven't updated the code yet.
7 %
8 % This function plots the risk of group 1 compared to group 2/3 for
9 % for many different values of two HRC's indicated by var1 and var2.
10 % The data for the groups are located in file1, file2, and file3. If
11 % file3 does not exist, the risk is defined as group 1 compared to
12 % group 2. If file3 exists, the risk is defined as the ratio of group 1
13 % compared to group 2 divided by group 3 compared to group 2. For more
14 % information about how we calculate the risk, please see the
15 % documentation for one_risk_matrix.m.
16 %
17 % Arguments:
18 % 1. preprocessing - a string indicating which preprocessing method
19 % to use
20 % 2. var1 - a number indicating the first HRC
21 % 3. var2 - a number indicating the second HRC
22 % 4. xmin - the minimum value for both HRCs
23 % 5. xmax - the maximum value for both HRCs
24 % 6. file1 - the file containing the half hours for group 1
25 % 7. file2 - the file containing the half hours for group 2
26 % 8. file3 - the file containing the half hours for group 3
27 %
28 % Preconditions:
29 % 1. Make sure the files one_risk_matrix.m, one_porb_matrix.m,
30 % hrc_indices.m, and bin_widths.m are in the working directory.
31 %
32 % Returns:
33 % 1. This function returns a 3D plot of the risks of group 1
34 % compared to group 2/3 for many different values of the two
35 % HRCs.
36 %
37 % First, we read in the data.
38 csv_dir = [pwd '/Data-Files/Dienstman_CSV-Files-' preprocessing '.PP//'];
39
40 half_hours_group_1 = dlmread([csv_dir file1], ',', 1, 0);
41 half_hours_group_2 = dlmread([csv_dir file2], ',', 1, 0);
42
43 if exist('file3', 'var')
44 organism = str2double(file1(15));
45 half_hours_group_2 = half_hours_group_2(half_hours_group_2(:,10) ~= organism, :);
46 half_hours_group_3 = dlmread([csv_dir file3], ',', 1, 0);
47 end
48
49 % Next, we calculate the risk for different HRC values.
50 if exist('file3', 'var')
51 [risk_matrix, points] = one_risk_matrix(preprocessing, [var1 var2], xmin, xmax,
    half_hours_group_1, half_hours_group_2, half_hours_group_3);
52 else
53 [risk_matrix, points] = one_risk_matrix(preprocessing, [var1 var2], xmin, xmax,
    half_hours_group_1, half_hours_group_2);
54 end
55
56 % Finally, we plot the risks. Note that we use a log scale. Thus, any
57 % risk over 0 indicates the HRC values are more likely to be in
58 % group 1. Any risk below 0 indicates the HRC values are more likely to

```

```

59 % be in group 2/3.
60 [~, axes_labels] = xlsread([csv_dir file1], '1:1');
61 [X, Y] = meshgrid(points(1,:), points(2,:));
62 surf(X, Y, risk_matrix)
63 xlabel(regexprep(axes_labels{var1}, '-1', ' '))
64 ylabel(regexprep(axes_labels{var2}, '-1', ' '))
65 zlabel('Elevated Risk')
66 zlim([log(.5) log(2)])
67 caxis([log(.5) log(2)])
68 end

```

B.8 Single Variable Logistic Figures

Listing B.26: multiple_univariate_logistic_figures.m

```

1 % Author: Evan Dienstman
2 % Last Update: 3/31/2017
3 % Email: eddienstman@email.wm.edu
4 % Note: Feel free to email me with questions! If something doesn't
5 % make sense, it might be because I haven't updated the code yet.
6 %
7 % This script makes multiple figures that plot the probability of
8 % illness for various values of each HRC. We first start off by making
9 % figures for on specific category of half hours. For example, we might
10 % only use the nonventilated half hours. Next, we make five figures for
11 % each of these half hour categories corresponding to the five HRC
12 % subcategories. Each figure then contains seven subplots corresponding
13 % to the seven HRCs. On each subplot, we graph the probability of
14 % illness calculated from single variable logistic regression and from
15 % bayesian methods. We also graph the univariate PDFs for the HRCs on
16 % the subplot to give a sense of where most of the data lies. For
17 % example, the PDF of variance from 0-1 is very small. Therefore,
18 % even though a variance between 0-1 indicates an increased
19 % probability of being ill, we also want to note that these values are
20 % very rare to begin with. For extremely rare values, we do not even
21 % plot the probability of illness. For more info about the figures,
22 % see the documentation of one_univariate_logistic_plot.m.
23 %
24 % Preconditions:
25 % 1. Make sure the files one_univariate_logistic_figure,
26 % one_risk_matrix.m, one_porb_matix.m, hrc_indices.m, and
27 % hrc_bin_widths.m are in the current working directory.
28 % 2. Make sure all the CSV files are in the appropriate directory.
29 % 3. Make sure to select the proper preprocessing method below.
30 %
31 % Returns:
32 % 1. This script returns multiple single variable logistic figures
33 % saved to the directory indicated in
34 % one_univariate_logistic_figure.m with the file name given
35 % in the call to one_univariate_logistic_figure.m.
36
37 clear
38 clc
39
40 % Here, we select the preprocessing method.
41 preprocessing = 'Abby';
42 % preprocessing = 'Doug';
43
44 % Next, we define the CSV files we want to use for the figures. Note
45 % that we don't consider and specific organism categories because the
46 % PDF and risk figures told us that the HRCs are not useful for
47 % distinguishing amongst different organisms.
48 file_list_1 = {'hrc_sick_org_all_vent_all.csv', 'hrc_healthy_org_all_vent_all.csv'};
49 file_list_2 = {'hrc_sick_org_all_vent_0.csv', 'hrc_healthy_org_all_vent_0.csv'};
50 file_list_3 = {'hrc_sick_org_all_vent_1.csv', 'hrc_healthy_org_all_vent_1.csv'};
51
52 % Finally, we call one_univariate_logistic_figure to make the
53 % figures for each group of CSV files.
54 one_univariate_logisitic_figure(preprocessing, file_list_1, 'hrc_all_org_all_vent_all')
55 one_univariate_logisitic_figure(preprocessing, file_list_2, 'hrc_all_org_all_vent_0')
56 one_univariate_logisitic_figure(preprocessing, file_list_3, 'hrc_all_org_all_vent_1')

```

Listing B.27: one_univariate_logistic_figure.m

```

1 function one_univariate_logisitic_figure(preprocessing, file_list, save_str)
2 % Author: Evan Dienstman
3 % Last Update: 3/31/2017
4 % Email: eddienstman@email.wm.edu
5 % Note: Feel free to email me with questions! If something doesn't
6 % make sense, it might be because I haven't updated the code yet.
7 %
8 % This script makes fives figures that plot the probability of
9 % illness for various values of each HRC using the CSV files in
10 % file_list. Each of the five figures corresponds to one of the five

```

```

11 % HRC subcategories (raw, 10th, 50th, 90th, and slope). Each figure
12 % then contains seven subplots corresponding to the seven HRCs. On each
13 % subplot, we graph the probability of illness calculated from single
14 % variable logistic regression and from bayesian methods. We also
15 % graph the univariate PDFs for the HRCs on the subplot to give a sense
16 % of where most of the data lies. For example, the PDF of variance
17 % from 0-1 is very small. Therefore, even though a variance between 0-1
18 % indicates an increased probability of being ill, we also want to note
19 % that these values are very rare to begin with. For extremely rare
20 % values, we do not even plot the probability of illness. For more info
21 % about the figures, see the comments below.
22 %
23 % Arguments:
24 % 1. preprocessing - the type of preprocessing used to create the
25 % CSV files
26 % 2. file_list - the list of CSV files used to make the figure
27 % 3. save_file - the name of the file the figure is saved to
28 %
29 % Preconditions:
30 % 1. Make sure the files one_risk_matrix.m, one_porb_matix.m,
31 % hrc_indices.m, and hrc_bin_widths.m are in the current working
32 % directory.
33 % 2. Make sure all the CSV files are in the appropriate directory.
34 % 3. This script will not overwrite any existing files with the
35 % same name. Change the variable save_str to a name that does
36 % not already exist, or delete the existing file with the same
37 % name before running this script.
38 %
39 % Returns:
40 % 1. This function returns five figures corresponding to the five
41 % HRC subcategories. The figures are saved to the directory
42 % indicated by the variable save_dir with the file name given
43 % by the variable save_file.
44 %
45 % Frist, load hrc_indices.mat and define some constants and variables
46 % used later.
47 load('hrc_indices.mat')
48 model = 'quadratic';
49 hrc_names = {'Asymmetry-1', 'Asymmetry-2', 'Asymmetry-Ratio', 'Decelerations', 'Mean-RR', '
    Sample-Entropy', 'Variance'};
50 hrc_types = {'', '_10', '_50', '_90', '_Slope'};
51 title_strs = {'Raw HRC', '10th Percentile', '50th Percentile', '90th Percentile', 'HRC Slope'};
52 x_axis_list = {[0 10], [0 10], [0 5], [0 50], [250 600], [0 1.5], [0 5]};
53 x_axis_list_slope = {[ -5 5], [ -5 5], [ -1 1], [ -10 10], [ -50 50], [ -1 1], [ -2 2]};
54 %
55 % Next, we read in the data from the CSV files.
56 csv_dir = [pwd '/Data-Files/Dienstman-CSV-Files-' preprocessing '_PP'];
57 csv_matrix_1 = dlmread([csv_dir '/' file_list{1}], ',', 1, 0);
58 csv_matrix_2 = dlmread([csv_dir '/' file_list{2}], ',', 1, 0);
59 csv_matrix = [csv_matrix_1; csv_matrix_2];
60 healthy_half_hours = csv_matrix(csv_matrix(:, health_status_index) == 0,:); %ok< *NASGU>
61 sick_half_hours = csv_matrix(csv_matrix(:, health_status_index) == 1,:);
62 p_sick = length(sick_half_hours)/length(csv_matrix);
63 %
64 % We now loop through all five HRC subcategories and create a figure
65 % for each one.
66 for i = 1:length(hrc_types)
67 %
68 % Here, we create the save directory and save file for the figure.
69 save_dir = [pwd '/Figure-Files/Dienstman-Logistic-Figures-' preprocessing '_PP'];
70 save_file = [save_dir '/' save_str hrc_types{i} '.png'];
71 %
72 % If the save_dir doesn't already exist, we make the save_dir here.
73 if ~exist(save_dir, 'dir')
74 mkdir(save_dir)
75 end
76 %
77 % If the save file already exists, we stop the program so we don't
78 % overwrite the file.
79 if exist(save_file, 'file')
80 disp('Error: A file already exists with the save file name. The program stopped because
    running the program would overwrite the existing file.')
81 continue
82 end
83 %
84 % Here, we create the figure, info box in the top left corner, and
85 % the annotation at the bottom of the figure.
86 probability_figure = figure('Position', [50,50,1600,900]);
87 set(probability_figure, 'color', 'w');
88 title_frame = uicontrol('style', 'frame');
89 if i >= 2 && i <= 4
90 baby_title_string = [save_str ' ' title_strs{i} char(10) char(10)];
91 else
92 baby_title_string = [' ' save_str ' ' title_strs{i} char(10) char(10)];
93 end
94 baby_title_string = [baby_title_string ...
95 ' Probability from Healthy (0) to Sick (1)' char(10) ...
96 ' Red Line = Bayesian Probability' char(10) ...
97 ' Blue Line = Logistic Probability' char(10) char(10) ...
98 ' Fraction from No Data (0) to All Data (1)' char(10) ...
99 ' Red Area = Sick Fraction' char(10) ...
100 ' Blue Area = Healthy Fraction']; %ok< *AGROW>
101 set(title_frame, 'Position', [210, 635, 340, 195], 'BackgroundColor', [0 0 0])
102 baby_title = uicontrol('style', 'text');

```

```

103 set(baby_title, 'String', baby_title_string, 'Position', [213, 638, 334, 189], 'FontSize', 13,
104     'BackgroundColor', [1 1 1])
105 axes_note_string = 'All red and blue areas are additive, i.e., no areas are hidden behind one
106     another.';
107 axes_note = uicontrol('style', 'text');
108 set(axes_note, 'String', axes_note_string, 'Position', [530, -20, 600, 50], 'FontSize', 12, '
109     'BackgroundColor', [1 1 1])
110
111 % Here, we iterate through all seven HRCs and create the
112 % corresponding probability plot for that HRC.
113 for j = 1:length(hrc_names)
114
115     % First, we specify the subplot.
116     subplot(3,3,j+1)
117     var_name = [hrc_names{j} hrc_types{i}];
118     var_index = [lower(var_name) '_index'];
119
120     % Next, we separate the data into the HRC values (data_vector)
121     % and health status (response_vector).
122     eval(['data_vector = csv_matrix(:, ' var_index ');'])
123     response_vector = csv_matrix(:, health_status_index);
124     [data_vector, sorted_indices] = sort(data_vector);
125     response_vector = response_vector(sorted_indices);
126
127     % Before we calculate the probabilities, we remove outliers
128     % from the data. For the Bayesian method, removing outliers
129     % will not affect the results because outliers will have very
130     % low probabilities. However, we want to remove outliers for
131     % the logistic probability because over fitting will
132     % considerably change the tails if we keep the outliers. For
133     % decelerations, we define the high_fence as 30 because the
134     % outlier method removes too much data. This procedure is
135     % strictly empirical and needs to be analyzed further.
136     q1 = quantile(data_vector, .25);
137     q3 = quantile(data_vector, .75);
138     IQR = q3 - q1;
139     low_fence = q1 - 1.5*IQR;
140     high_fence = q3 + 1.5*IQR;
141
142     if strcmp(hrc_names{j}, 'Decelerations')
143         data_vector(data_vector < low_fence | data_vector > 30) = NaN;
144         response_vector(data_vector < low_fence | data_vector > 30) = NaN;
145     else
146         data_vector(data_vector < low_fence | data_vector > high_fence) = NaN;
147         response_vector(data_vector < low_fence | data_vector > high_fence) = NaN;
148     end
149
150     % Now we calculate the probability of illness using logistic
151     % regression. We use the Matlab function fitglm to calculate the
152     % probability. For more information, please see the Matlab
153     % documentation of this function.
154     fit = fitglm(data_vector, response_vector, model, 'distribution', 'binomial', 'VarNames',
155         {var_name, 'Health_Status'});
156     log_prob = fit.Fitted.Probability;
157
158     % Now we calculate the probability of illness using the Bayesian
159     % probability. We use the function one_risk_matrix to
160     % calculate the probability. For more information, please see
161     % the documentation of this function.
162     if strcmp(title_strs{i}, 'HRC Slope')
163         eval(['risk, x_points, p_sigs_given_sick, p_sigs_given_healthy] = one_risk_matrix(
164             preprocessing, ' var_index ', x_axis_list_slope{j}(1), x_axis_list_slope{j}(2),
165             sick_half_hours, healthy_half_hours);'])
166     else
167         eval(['risk, x_points, p_sigs_given_sick, p_sigs_given_healthy] = one_risk_matrix(
168             preprocessing, ' var_index ', x_axis_list{j}(1), x_axis_list{j}(2), sick_half_hours
169             , healthy_half_hours);'])
170     end
171     bayesian_prob = p_sick .* exp(risk);
172
173     % In each subplot, we also include a PDF for the sick and
174     % healthy half hours. With this information, we can visualize
175     % how much data we have at each HRC value. We want this
176     % information because we are more confident in probabilities at
177     % values with more data.
178     hold on
179     yyaxis right
180     myarea = area(x_points, [p_sigs_given_sick + p_sigs_given_healthy, p_sigs_given_healthy],
181         'FaceAlpha', 0.25, 'LineStyle', 'none');
182     myarea(1).FaceColor = 'B';
183     myarea(2).FaceColor = 'R';
184     ylim([0 1])
185     ylabel('Fraction of Data')
186     yyaxis left
187     plot(data_vector, log_prob, 'LineWidth', 2, 'Color', 'B')
188     plot(x_points, bayesian_prob, '-', 'LineWidth', 2, 'Color', 'R')
189     ylabel('Probability')
190     ylim([0 0.025])
191
192     % Lastly, we set and label the axes for the subplot.
193     if strcmp(hrc_names{j}, 'Mean-RR') || strcmp(title_strs{i}, 'HRC Slope')
194         xmin = min([data_vector; x_points]);
195     else
196         xmin = 0;
197

```

```

188     end
189
190     xmax = max([data_vector; x_points]);
191     xlim([xmin xmax])
192     xlabel(regexprep(hrc_names{j}, '-1', ' '))
193     set(gca, 'FontSize', 14)
194     hold off
195 end
196
197 % Finally, we save the figure.
198 print(probability_figure, save_file, '-dpng')
199 end
200 end

```

B.9 HeRO Score Figures

Listing B.28: `multiple_hero_score_figures.m`

```

1 % Author: Evan Dienstman
2 % Last Update: 4/7/2017
3 % Email: eddienstman@email.wm.edu
4 % Note: Feel free to email me with questions! If something doesn't
5 % make sense, it might be because I haven't updated the code yet.
6 %
7 % This script creates a hero score figure for every septic events.
8 % Each hero score figure shows six hero scores at each half hour of
9 % the Dienstman_results file associated with the event in question.
10 % The six hero scores are Dienstman_all, Dienstman_vent,
11 % Dienstman_nonvent, Hrch, Hrcg, and Hero. Each hero score is calculated
12 % using the half hours in the Dienstman_results files. Starting from
13 % the half hour in question, the hero score looks at the prior 24 half
14 % hours when considering what heart rate characteristics (HRCs) to use
15 % in the score. Thus, the first hero score in the figure looks at half
16 % hours 1-24 in the Dienstman_results file, the second hero score looks
17 % at half hours 2-25, and so on. The verticle black line in the figure
18 % represents the time of the event. For more information about how to
19 % calculate the hero score, see the documentations for
20 % one_hero_score_figure. This function then creates a figure containg
21 % the average hero scores calculated from all the individual figures.
22 %
23 % Preconditions:
24 % 1. Make sure the directories and file names used in the script
25 % are the right ones for the computer you are using.
26 % 2. Make sure the files event_matrix.mat, Doug_coeffs.mat,
27 % Dienstman_coeffs.mat, prctile1.m, logistic.m, and
28 % one_hero_score_figure.m are in the working directory.
29 % 3. This function will not overwrite the average figure that
30 % already exists. Delete the old figure or change the name of
31 % save_file below.
32 %
33 % Returns:
34 % 1. This script returns a hero score figures for every septic event
35 % and an average figure for all the events.
36
37 clear
38 clc
39
40 % Change the preprocessing to the one you want to use.
41 preprocessing = 'Abby';
42 % preprocessing = 'Doug';
43
44 % First, we create the save_dir and save_file.
45 save_dir = [pwd '/Figure_Files/Dienstman_Hero_Scores_' preprocessing '_PP/Average'];
46 save_file = [save_dir '/average_hero_figure.fig'];
47
48 if ~exist(save_dir, 'dir')
49     mkdir(save_dir)
50 end
51
52 % If the save file already exists, we stop the program so we don't
53 % overwrite the file.
54 if exist(save_file, 'file')
55     disp(['Error: File ' save_file ' already exists. Did not execute because program would alter
56         existing file.'])
57     return
58 end
59
60 % Here, we create the container that let's us translate between the
61 % site codes and the site numbers. The site refers to the hospitals of
62 % the babies.
63 site_map_keys = {11, 13, 15, 23, 24, 26, 27, 30};
64 site_map_values = {'UVA', '0d', '0f', '17', '18', '1a', '1b', '1e'};
65 site_map = containers.Map(site_map_keys, site_map_values);
66 plot_str = 'yes';
67
68 % Here, we load in the variable event_matrix and hrc_indices. The
69 % variable event_info contains the site, ID, time, organism, and

```

```

69 % other demographic info for each spetic event. The variable
70 % hrc_indices contains the column indices of the CSV files for
71 % each HRC.
72 load('event_matrix.mat')
73 load('hrc_indices.mat')
74
75 % Next, we define some variables to use later.
76 index_vector = -335:144;
77 mean_dienstman_hero_vector_all = zeros(1,480);
78 mean_dienstman_hero_vector_vent = zeros(1,480);
79 mean_dienstman_hero_vector_nonvent = zeros(1,480);
80 mean_hrch_vector = zeros(1,480);
81 mean_hrcg_vector = zeros(1,480);
82 mean_hero_score_vector = zeros(1,480);
83 all_count_vector = zeros(1,480);
84 vent_count_vector = zeros(1,480);
85 nonvent_count_vector = zeros(1,480);
86 count = 1;
87
88 % Next, we interate through every index in rand_indices and create
89 % a figure for the event corresponding to that index in event_matrix.
90 for rand_index = 1:length(event_matrix)
91     id = event_matrix(rand_index,1);
92     site_num = event_matrix(rand_index,2);
93     site = site_map(site_num);
94     event_time = event_matrix(rand_index,7);
95
96     save_dir = [pwd '/Figure_Files/Dienstman_Hero_Scores_' preprocessing '_PP/' site];
97     Dienstman_dir = [pwd '/Data_Files/Dienstman_Results_' preprocessing '_PP/' site];
98     Dienstman_file = [Dienstman_dir '//Dienstman_hrc_results_' site '_' num2str(id) '.mat'];
99
100 % If the save directory doesn't exist, we make it here.
101 if ~exist(save_dir, 'dir')
102     mkdir(save_dir)
103 end
104
105 % Lastly, we call one_hero_score_figure with the information
106 % for this event to create the hero score figure.
107 [one_dienstman_hero_vector_all, one_dienstman_hero_vector_vent,
    one_dienstman_hero_vector_nonvent, one_hrch_vector, one_hrcg_vector, time_vector] =
    one_hero_score_figure(id, site_num, event_time, save_dir, Dienstman_file, plot_str);
108
109 % After we create the individual figure, we update the average
110 % vectors.
111 for i = 1:480
112     start_time = (event_time - 7) + i/48;
113     end_time = (event_time - 7) + (i+1)/48;
114     index = find(time_vector > start_time & time_vector < end_time);
115
116     if length(index) == 1 && ~isnan(one_dienstman_hero_vector_all(index))
117         mean_dienstman_hero_vector_all(i) = nansum([mean_dienstman_hero_vector_all(i),
118             one_dienstman_hero_vector_all(index)]);
119         mean_dienstman_hero_vector_vent(i) = nansum([mean_dienstman_hero_vector_vent(i),
120             one_dienstman_hero_vector_vent(index)]);
121         mean_dienstman_hero_vector_nonvent(i) = nansum([mean_dienstman_hero_vector_nonvent(i),
122             one_dienstman_hero_vector_nonvent(index)]);
123         mean_hrch_vector(i) = nansum([mean_hrch_vector(i), one_hrch_vector(index)]);
124         mean_hrcg_vector(i) = nansum([mean_hrcg_vector(i), one_hrcg_vector(index)]);
125         mean_hero_score_vector(i) = nansum([mean_hero_score_vector(i), max(one_hrch_vector(
126             index), one_hrcg_vector(index))]);
127         all_count_vector(i) = all_count_vector(i) + 1;
128
129         if ~isnan(one_dienstman_hero_vector_vent(index))
130             vent_count_vector(i) = vent_count_vector(i) + 1;
131         end
132
133         if ~isnan(one_dienstman_hero_vector_nonvent(index))
134             nonvent_count_vector(i) = nonvent_count_vector(i) + 1;
135         end
136     end
137 end
138 count = count + 1;
139 end
140
141 % Here, we calculate the average vectors.
142 mean_dienstman_hero_vector_all = mean_dienstman_hero_vector_all./all_count_vector;
143 mean_dienstman_hero_vector_vent = mean_dienstman_hero_vector_vent./vent_count_vector;
144 mean_dienstman_hero_vector_nonvent = mean_dienstman_hero_vector_nonvent./nonvent_count_vector;
145 mean_hrch_vector = mean_hrch_vector./all_count_vector;
146 mean_hrcg_vector = mean_hrcg_vector./all_count_vector;
147 mean_hero_score_vector = mean_hero_score_vector./all_count_vector;
148
149 % Finally, we plot the average vectors.
150 figure_handle = figure('Position', [50,50,1600,900]);
151 set(figure_handle, 'color', 'w');
152 hold on
153 plot(index_vector, mean_dienstman_hero_vector_all, 'LineWidth', 2, 'Color', 'R')
154 plot(index_vector, mean_dienstman_hero_vector_vent, 'LineWidth', 2, 'Color', [1 1 0])
155 plot(index_vector, mean_dienstman_hero_vector_nonvent, 'LineWidth', 2, 'Color', [0 1 1])
156 plot(index_vector, mean_hrch_vector, 'LineWidth', 2, 'Color', 'B')
157 plot(index_vector, mean_hrcg_vector, 'LineWidth', 2, 'Color', 'G')
158 plot(index_vector, mean_hero_score_vector, 'LineWidth', 2, 'Color', [1 0 1])
159 line([0 0], ylim, 'Color', [0,0,0], 'LineWidth', 3);

```

```

157 legend('Dienstman Hero Score', 'Dienstman Hero Score Vent', 'Dienstman Hero Score Nonvent', 'Hrch
      Sore', 'Hrcg Score', 'Hero Score')
158 title(['Average HeRO Score: Mean Babies Used Per Half Hour - ' num2str(mean(all_count_vector))'], '
      FontSize', 24)
159 xlabel('Time of HeRO Score', 'FontSize', 16)
160 ylabel('Hero Score', 'FontSize', 16)
161 set(gca, 'fontsize', 16)
162 hold off
163
164 % Lastly, we save the average figure.
165 hgsave(figure_handle, save_file, '-v7.3');

```

Listing B.29: one_hero_score_figure.m

```

1 function [dienstman_hero_vector_all, dienstman_hero_vector_vent, dienstman_hero_vector_nonvent,
      hrch_vector, hrcg_vector, plot_time_vector] = one_hero_score_figure(id, site_num, event_time,
      save_dir, Dienstman_file, plot_str)
2 % Author: Evan Dienstman
3 % Last Update: 4/7/2017
4 % Email: eddienstman@email.wm.edu
5 % Note: Feel free to email me with questions! If something doesn't
6 % make sense, it might be because I haven't updated the code yet.
7 %
8 % This function creates a HeRO score figure for one septic event. We
9 % save the figure to the directory Dienstman_Hero_Scores. Each figure
10 % contains the multiple HeRO scores calculated at each half hour. The
11 % HeRO scores include Dienstman_Hero_All, Dienstman_Hero_Vent,
12 % Dienstman_Hero_Nonvent, Hero, Hrch, and Hrcg. For more info on these
13 % six HeRO scores, please see the code below. We calculate the HeRO
14 % scores using the heart rate characteristics (HRCs) for each half hour.
15 % The HRCs for each half hour are stored in the Dienstman_results files.
16 % The HeRO score looks over a window of half hours in the past given by
17 % the variable half_hour_window_length found in the code below. The
18 % HeRO score then takes a certain percentile value for each HRC to
19 % calculate the HeRO score. Since we are looking at a window into the
20 % past, the time of the first HeRO score in the figure is the time of
21 % the half hour in the Dienstman_file with an index equal to the number
22 % half_hour_window_length. Note that the percentiles of the HRCs have
23 % already been calculated in the Dienstman_result files. Thus, we do
24 % not need to calculate them in this file. The black verticle line
25 % represents the time of the event.
26 %
27 % Arguments:
28 % 1. id - the id number of the baby
29 % 2. site - the site of the baby
30 % 3. event_time - the time of the septic event for the baby
31 % 4. save_dir - the directory where we will save the figure
32 % 5. Dienstman_file - the Dienstman_file that corresponds to the
33 % event that we will calculate the HeRO scores for
34 %
35 % Preconditions:
36 % 1. Make sure the directories and file names used in the scripts
37 % are the right ones for the computer you are using.
38 % 2. Make sure the files Dienstman_coeffs.mat, Doug_coeffs.mat,
39 % prctile1.m, and logistic.m, are in the working directory.
40 % 3. Make sure the variable half_hour_window_length matches the
41 % time window you want the HeRO score to calculate.
42 % 4. This function will not overwrite any figure that already
43 % exists. Delete the old figure or change the name of save_file
44 % below.
45 %
46 % Returns:
47 % 1. dienstman_hero_vector_all - the vector containing the
48 % Dienstman_Hero_All HeRO score
49 % 2. dienstman_hero_vector_vent - the vector containing the
50 % Dienstman_Hero_Vent HeRO score
51 % 3. dienstman_hero_vector_nonvent - the vector containing the
52 % Dienstman_Hero_Nonvent HeRO score
53 % 4. hrch_vector - the vector containing the Hrch HeRO score
54 % 5. hrcg_vector - the vector containing the Hrcg HeRO score
55 % 6. plot_time_vector - the vector containing the start times
56 % of the half hours
57 % 7. This function also returns a figure with all six HeRO scores at
58 % every half hour of the file Dienstman_file.
59 %
60 % First, we define some variables that we will use later.
61 site_map_keys = {11, 13, 15, 23, 24, 26, 27, 30};
62 site_map_values = {'UVA', '0d', '0f', '17', '18', '1a', '1b', '1e'};
63 site_map = containers.Map(site_map_keys, site_map_values);
64 site = site_map(site_num);
65 save_file = [save_dir '/' hero_figure_' site '_' num2str(id) '_' num2str(round(event_time)) '.fig'
      ];
66
67 % If the save file already exists, we stop the program so we don't
68 % overwrite the file.
69 if exist(save_file, 'file')
70     disp(['Error: File ' save_file ' already exists. Did not execute because program would alter
      existing file.'])
71     hrch_vector = NaN;
72     hrcg_vector = NaN;
73     dienstman_hero_vector_all = NaN;
74     dienstman_hero_vector_vent = NaN;

```



```

75     dienstman_hero_vector_nonvent = NaN;
76     plot_time_vector = NaN;
77     return
78 end
79
80 % We then load in the Dienstman file containing the half hours needed
81 % to calculate the hero scores for this event. If the Dienstman-file
82 % doesn't exist, we return to the calling function.
83 if exist(Dienstman_file, 'file')
84     load_variable = load(Dienstman_file);
85     struct_name = fieldnames(load_variable);
86     struct_name = struct_name(1);
87     eval_str = ['Dienstman_struct = load_variable.' char(struct_name) ';''];
88     eval(eval_str)
89 else
90     hrch_vector = NaN;
91     hrcg_vector = NaN;
92     dienstman_hero_vector_all = NaN;
93     dienstman_hero_vector_vent = NaN;
94     dienstman_hero_vector_nonvent = NaN;
95     plot_time_vector = NaN;
96     return
97 end
98
99 % Here, we load in variables needed to calculate the hero score. Each
100 % one of the variables is a coefficient in the hero score model that
101 % has been optimized to give the best results.
102 load Doug_coeffs cg ch u0
103 load('Dienstman_coeffs_all.mat')
104 load('Dienstman_coeffs_vent.mat')
105 load('Dienstman_coeffs_nonvent.mat')
106 load('vent_matrix.mat')
107
108 % Change this number if you want the HeRO score to encompass more or
109 % less half hours in the past. Currently, it is set to 24 to encompass
110 % 24 half hours in the past.
111 half_hour_window = 24;
112
113 % Using the Dienstman_file we loaded, we create the time_vector each
114 % hero score will be calculated at. Note that we start the time
115 % starting at half_hour_window.length because the first score will
116 % look at that number of half hours in the past. We then preallocate
117 % the vectors below that we will plot later.
118 time_vector = [Dienstman_struct(:).Start_Time]; %%ok< *NODEF>
119 plot_time_vector = time_vector(half_hour_window:end);
120 num_of_hero_scores = length(time_vector) - half_hour_window + 1;
121 hero_score_vector = zeros(1, num_of_hero_scores);
122 hrch_vector = zeros(1, num_of_hero_scores);
123 hrcg_vector = zeros(1, num_of_hero_scores);
124 dienstman_hero_vector_all = zeros(1, num_of_hero_scores);
125 dienstman_hero_vector_vent = zeros(1, num_of_hero_scores).*NaN;
126 dienstman_hero_vector_nonvent = zeros(1, num_of_hero_scores).*NaN;
127 vent_indices = find(vent_matrix(:,1) == site_num & vent_matrix(:,2) == id); %%ok< *NODEF>
128 baby_vent_info = vent_matrix(vent_indices,3:4); %%ok< *FNDSB>
129
130 % Now we iterate through every half hour, calculating the hero score
131 % each time. Again, we start at the half hour with index
132 % half_hour_window.length because the hero score looks at that
133 % number of half hours in the past.
134 for i = half_hour_window:(num_of_hero_scores + half_hour_window - 1)
135     hrc_entry = Dienstman_struct(i);
136     indices = find(time_vector <= hrc_entry.Start_Time & time_vector >= (hrc_entry.Start_Time -
137         half_hour_window/48));
138
139     if isempty(hrc_entry.Start_Time) && isempty(indices)
140         hrc_struct = Dienstman_struct(indices);
141
142         % First, we calculate the Hrch, Hrcg, and Hero scores. The
143         % components of these HeRO scores are contained in the vector
144         % hero_score_nums_1. These HeRO scores look at different
145         % percentiles of HRCs over the half_hour_window. The
146         % coefficients for these HeRO scores are contained in
147         % Doug_coeffs.mat
148         hero_score_nums_1 = zeros(1,6);
149         hero_score_nums_1(1) = hrc_entry.Variance_10;
150         hero_score_nums_1(2) = hrc_entry.Sample_Entropy_10;
151         hero_score_nums_1(3) = hrc_entry.Asymmetry_1.50;
152         hero_score_nums_1(4) = hrc_entry.Asymmetry_2.50;
153         hero_score_nums_1(5) = prctile1([hrc_struct(:).Asymmetry_2] - [hrc_struct(:).Asymmetry_1],
154             50);
155         hero_score_nums_1(6) = hrc_entry.Variance_50;
156         hrch = logistic(ch, hero_score_nums_1([1 5])) / u0;
157         hrcg = logistic(cg, hero_score_nums_1([6 2 3 4])) / u0;
158         hero_score = max(hrch, hrcg);
159
160         % Next, we determine if the half hour is ventilated.
161         if isempty(find(hrc_entry.Start_Time > baby_vent_info(:,1) & hrc_entry.Start_Time <
162             baby_vent_info(:,2), 1))
163             ventilated = 'yes';
164         else
165             ventilated = 'no';
166         end
167
168         % Now we calculate the Dienstman_Hero_Vent or

```

```

166 % Dienstman_Hero_Nonvent HeRO score depending on the
167 % ventilation status of the half hour as well as the
168 % Dienstman_Hero_All score.
169 variables = {'Asymmetry_1_10', 'Asymmetry_1_Slope', ...
170 'Asymmetry_2_90', 'Asymmetry_2_Slope', ...
171 'Asymmetry_Ratio_50', 'Asymmetry_Ratio_Slope', ...
172 'Decelerations_90', ...
173 'Sample_Entropy_10', 'Sample_Entropy_Slope', ...
174 'Variance_10', 'Variance_90', 'Variance_Slope'};
175 N = length(variables);
176 num_of_coefs = sum(1:N) + N;
177 hero_score_nums_2 = zeros(1,N);
178 hero_score_nums_3 = zeros(1,N);
179
180 if strcmp(ventilated, 'no')
181     coeffs2 = nonvent_coefs;
182 else
183     coeffs2 = vent_coefs;
184 end
185
186 % Here, we set the components of the Dienstman_Hero_Vent/Nonvent
187 % and Dienstman_Hero_All HeRO scores.
188 for j = 1:N
189     eval(['hero_score_nums_2(j) = hrc_entry.' variables{j} ';''])
190     eval(['hero_score_nums_3(j) = hrc_entry.' variables{j} ';''])
191 end
192
193 % Since we have a quadratic model, we now calculate all the
194 % cross products from the HRC components. We save these new
195 % components in Dienstman_hero_score_nums_2. Note that the
196 % order we save these components is important because they
197 % must match the coefficients in Dienstman_coefs_vent/nonvent
198 % and Dienstman_coefs_all.
199 Dienstman_hero_score_nums_2 = zeros(1, num_of_coefs);
200 Dienstman_hero_score_nums_3 = zeros(1, num_of_coefs);
201 count = 1;
202
203 for j = 1:N
204     Dienstman_hero_score_nums_2(count) = hero_score_nums_2(j);
205     Dienstman_hero_score_nums_3(count) = hero_score_nums_3(j);
206     count = count + 1;
207 end
208
209 for j = 1:N
210     for k = j+1:N
211         Dienstman_hero_score_nums_2(count) = hero_score_nums_2(j) * hero_score_nums_2(k);
212         Dienstman_hero_score_nums_3(count) = hero_score_nums_3(j) * hero_score_nums_3(k);
213         count = count + 1;
214     end
215 end
216
217 for j = 1:N
218     Dienstman_hero_score_nums_2(count) = hero_score_nums_2(j)^2;
219     Dienstman_hero_score_nums_3(count) = hero_score_nums_3(j)^2;
220     count = count + 1;
221 end
222
223 % Finally, we calculate the hero score for
224 % Dienstman_Hero_Vent/Nonvent and Dienstman_Hero_all.
225 if sum(isnan(Dienstman_hero_score_nums_2)) == 0 && sum(isnan(Dienstman_hero_score_nums_3))
    == 0
226     Dienstman_hero_score_2 = logistic(coeffs2, Dienstman_hero_score_nums_2) / u0;
227     Dienstman_hero_score_3 = logistic(all_coefs, Dienstman_hero_score_nums_3) / u0;
228 else
229     Dienstman_hero_score_2 = NaN;
230     Dienstman_hero_score_3 = NaN;
231 end
232
233 % We then save the all the HeRO scores for plotting later.
234 hero_score_vector(i-half_hour_window+1) = hero_score;
235 hrch_vector(i-half_hour_window+1) = hrch;
236 hrcg_vector(i-half_hour_window+1) = hrcg;
237 dienstman_hero_vector_all(i-half_hour_window+1) = Dienstman_hero_score_3;
238
239 if strcmp(ventilated, 'no')
240     dienstman_hero_vector_nonvent(i-half_hour_window+1) = Dienstman_hero_score_2;
241 else
242     dienstman_hero_vector_vent(i-half_hour_window+1) = Dienstman_hero_score_2;
243 end
244 end
245 end
246
247 % Finally, we plot the HeRO scores at the time of each half hour.
248 if strcmp(plot_str, 'yes')
249     figure_handle = figure('Position', [50,50,1600,900]);
250     set(figure_handle, 'color', 'w');
251     hold on
252     plot_step = ceil(num_of_hero_scores * .01);
253     plot(plot_time_vector, dienstman_hero_vector_all, 'LineWidth', 4, 'Color', [.5 .5 .5])
254     plot(plot_time_vector, dienstman_hero_vector_vent, 'LineWidth', 4, 'Color', [135/255 206/255
255     250/255])
255     plot(plot_time_vector, dienstman_hero_vector_nonvent, 'LineWidth', 4, 'Color', [0 0 205/255])
256     plot(plot_time_vector, hero_score_vector, 'LineWidth', 4, 'Color', 'G')
257     plot(plot_time_vector(1:plot_step:end), hrch_vector(1:plot_step:end), '+', 'MarkerSize', 10)

```

```

258 plot(plot_time_vector(1:plot_step:end), hrcg_vector(1:plot_step:end), 'o', 'MarkerSize', 10)
259 line([event_time event_time], ylim, 'Color', [0,0,0], 'LineWidth', 3);
260 legend('Dienstman Hero', 'Dienstman Vent Hero', 'Dienstman Nonvent Hero', 'Legacy Hero', 'Hrch
      Sore', 'Hrcg Score')
261 title(['Moving HeRO Score (ID: ' num2str(id) ', Site: ' site ')'], 'FontSize', 24)
262 xlabel('Time of HeRO Score (Days)', 'FontSize', 16)
263 ylabel('HeRO Score Value', 'FontSize', 16)
264 set(gca, 'fontsize', 16)
265 hold off
266 hgsave(figure_handle, save_file, '-v7.3');
267 end
268 end

```

Bibliography

- [1] M.-K. ENCYCLOPEDIA, *electrocardiogram*, November 2003.
- [2] A. A. FLOWER, J. R. MOORMAN, D. E. LAKE, AND J. B. DELOS, *Periodic heart rate decelerations in premature infants*, *Experimental Biology and Medicine*, 235 (2010), pp. 531–538.
- [3] D. FREEDMAN AND P. DIACONIS, *On the histogram as a density estimator:l2 theory*, *Wahrscheinlichkeitstheorie verw Gebiete*, 57 (1980), pp. 453–476.
- [4] M. P. E. A. GRIFFIN, *Abnormal heart rate characteristics preceding neonatal sepsis and sepsis-like illness*, *Pediatric Research*, 53 (2003), pp. 920–926.
- [5] B. P. E. A. KOVATCHEV, *Sample asymmetry analysis of heart rate characteristics with application to neonatal sepsis and systemic inflammatory response syndrome*, *Pediatric Research*, 54 (2003), pp. 892–898.
- [6] J. R. E. A. MOORMAN, *Mortality reduction by heart rate characteristic monitoring in very low birth weight neonates: a randomized trial*, *The Journal of pediatrics*, 159 (2011), pp. 900–906.
- [7] J. S. RICHMAN AND J. R. MOORMAN, *Physiological time-series analysis using approximate entropy and sample entropy*, *American Journal of Physiology-Heart and Circulatory Physiology*, 278 (2000), pp. H2039–H2049.
- [8] M. C. STAFF, *Diseases and conditions: Sepsis*, July 2014.
- [9] A. THOMPSON, *The electrocardiogram - part v*, 2010.
- [10] WIKIPEDIA, *Kernel density estimation*, March 2017.
- [11] —, *Logistic regression*, April 2017.