# Certification

I declare that this thesis was written by me under the guidance and counsel of my supervisors.

......................................................... Date..........................
Abduljewad Nuru Wele             Student

We certify that this is the true thesis report written by **Abduljewad Nuru Wele** under our supervision and we thus permit its presentation for assessment.

......................................................... Date..........................
Dr. Edmund NJERU NJAGI             Supervisor

......................................................... Date..........................
Prof. Dr. Geert MOLENBERGHS    Co-Supervisor

......................................................... Date..........................
Prof. Dr. Paul DENDALE             Co-Supervisor

# Dedication

I dedicate this thesis

to my lovely family who sacrified their life opportunities to educate me and provided me both financial and moral support

And

to Vlaamse InterUniversitaire Raad (VLIR) for giving me the chance to study MSc. in biostatistics

# Acknowledgements

First and foremost, next to God for his unlimited blessing, I would like to thank my supervisors, Dr. Edmund Njagi, Prof. Dr. Geert Molenberghs and Prof. Dr. Paul Dendale for their guidance and constructive suggestions. This work might not be fruitful without the professional advice of my friendly supervisors. It has been a pleasure to work with Dr. Edmund. He provided me with his professional remarks, repeatitive checking with patience and constructive comments. Thanks for all the assistance and comments.

My gratitude will extend to Vlaamse InterUniversitaire Raad (VLIR). VLIR is the first organization that gave me the acadamic opportunity to enhance my skills in the field of statistics. In addition to financial support, VLIR enabled me to create professional network with different researchers. I also want to thank all professors and staff members at censtat who helped me in the last two years .

Special mention goes to my classmates, Kurnia Wahyudi (MD), Negera Wakgari and Forsi Nwebim Boeyeo for thier special suggestions and constructive comments while preparing the work.

Most of all, I wish to express my love and gratitude to my beloved families for their unlimited advice and endless love throughout my life.

<div style="text-align: right">Abduljewad Nuru Wele</div>

September 10, 2014

Diepenbeek, Belgium

# Abstract

Telemonitoring in chronic heart failure management aided in self monitoring of day to day patients biomarkers such as blood pressure, weight and heart rate. The clinicians in the health center use these measurements to enhance timely intervention and thereby reduce rehospitalization. We study the post discharge profiles of the different biomarkers measurements availed to the health center via telemonitoring using flexible statistical model. Analysis is based on the use of generalized additve mixed model (GAMM) using p-splines. The flexible additive linear predictor which is based on smooth functions is used to capture the rather complex functional relation between the biomarkers and time. We include parametric random effects to account for extra sources of variability. The post discharge profiles are also studied to investigate whether there exists differences in evolutions based on the baseline covariates: sex, left ventricular ejection fraction (LVEF), heart rhythm, New York Heart Association (NYHA) functional classification of the heart failure and age. The GAMM provided a good deal of flexibility in describing the overall post discharge profiles of these biomarkers and it also better described the evolution of biomarkers across levels of covariates. The clinicians might use these approach to predict patient conditions and rehospitalization.

**Key words**: chronic heart failure (CHF), generalized additive mixed models (GAMM), p-splines, telemonitoring

# Contents

# List of Tables

# List of Figures

# List of Abbreviations

| | |
|---|---|
| AIC | Aikake Information Criterion |
| AM | Additive Models |
| AMM | Additive Mixed Models |
| AR(1) | Auto Regressive order one |
| b.p.m | beats per minute |
| CHF | Chronic Heart Failure |
| d.f | Degrees of Freedom |
| DBP | Diastolic Blood Pressure |
| EDF. | Effective Degrees of Freedom |
| $E_P$ | Effective number of parameters |
| GAM | Generalized Additive Models |
| GAMM | Generalized Additive Mixed Models |
| GLMM | Generalized Linear Mixed Models |
| GCV | Generalized Cross validation |
| OCV | Ordinary cross validation |
| HR | Heart rate |
| RLRT | Restricted Loglikelihood Ratio Test |
| LVEF | Left Ventricular Ejection Fraction |
| ML | Maximum Likelihood |
| MSEP | Mean Square Error Prediction |
| NYHA | New York Heart Association class |
| NTproBNP | N-terminal pro brain natriuretic peptide |
| P-IRLS | Penalized Iterative Re-weighted Least Squares |
| SBP | Systolic Blood Pressure |
| TM | Telemonitoring |
| $\chi^2_{(k)}$ | Chi-square distribution with k degree of freedom |
| 95%CI | 95 percent Confidence Interval |

# 1   Introduction

## 1.1   Telemonitoring in Chronic Heart Failure Management

Heart failure is related with substantial morbidity, mortality, and healthcare costs (Ho *et al.*, 1993 and Bui and Fonarow, 2012). Patients are often followed closely in heart failure clinics but admittance to a cardiology department is frequently needed, and the 30-day readmission rate is high (Jencks *et al.*, 2009). The clinical course of heart failure is characterized by recurrent hospitalizations due to fluid overload and/or worsening of renal function (Dendale *et al.*, 2011). A regular adjustment of treatment of CHF patients is needed to lower morbidity, mortality, and healthcare costs. The European Guidelines on heart failure stress that education about medication adherence, early warning signs of impending decompensation, support of self-care behaviour, and optimization of pharmacological and device therapy are the main aims of long term care of heart failure (Dickstein *et al.*, 2008).

Telehealthcare (or telemonitoring) may be a solution, contributing a way to identify and monitor subclinical congestion. The earlier identification and treatment of congestion together with improved coordination of care may prevent hospitalization. The potential use of telemonitoring to help in transmitting important parameters is acknowledged, and the important players in the chronic management of heart failure patients are the patient, the primary care physician, and the heart failure management team (Dickstein *et al.*, 2008). The classical approach used in heart failure teams, with regular telephone contacts is however, labour-intensive. Therefore, the question arises of whether a close collaboration between general practitioners (GPs) and the heart failure nurse and/or cardiologist, facilitated by modern communication technology allowing day-to-day follow-up of body weight, blood pressure, and heart rate, may result in an improved clinical outcome (Dendale *et al.*, 2011).

The telemonitoring machine provides not only measurement alerts but also a graph of the evolution of the biomarkers profile (Dendale *et al.*, 2011). Such profiles are useful in visualizing the long term trends of a biomarker for a particular patient. However, sometimes the observed longitudinal profiles shows complicated trends which cannot be easily described. Hence, some

techniques have to be adopted in order to capture such trajectories in an appropriate way. Having the longitudinal setting, one can use linear mixed models to analyse such measurements and quantify the evolution of the biomarkers. As stated in Fitzmaurice *et al.*, (2009), by measuring study participants repeatedly through time, longitudinal studies allow the direct study of temporal changes within individuals and the factors that influence the change. This can be analyzed using some variant of the mixed effects models which allow modeling and analysis of between and within individual variation (Pinheiro and Bates, 2000). However, there is a wide variety of challenges that arise in analyzing longitudinal data. Although such parametric models enjoy simplicity, they have suffered from inflexibility in modeling complicated relationships between the response and covariates in various longitudinal studies such as growth curves (Verbeke and Molenberghs, 2000). To circumvent the inflexibility of linear mixed models, a wide variety of models have been developed recently (Fitzmaurice *et al.*, 2009). Generalized additive models (GAMs) are among the alternatives. GAMs are flexible statistical models that captures the complex functional relationships of the outcomes and covariates in a smoothed manner rather than constant fashion (Wood, 2006).

In this project, the generalized additive mixed models are the primary method of analysis to describe longitudinal postdischarge profiles of blood pressure, weight and heart rate by capturing the functional relationships of these biomarkers and time.

## 1.2 Objectives

This project has two main objectives: (1) to study the postdischarge longitudianl profiles of the four biomarkers (diastolic, systolic, weight and heart rate) using flexible statistical techniques in such a way that one can easily see the long term trends. (2) to investigate whether the evolution of the profiles are influenced by the patients baseline characteristics (age, sex, LVEF, heart rhythm, NYHA and NTproBNP). Analyses will be conducted using package 'mgcv' in R.

## 2 Data description

### 2.1 The Chronic Heart Failure

The data used in this study is obtained from a study conducted in Belgium between 2008 and 2010 and whose aim was to study whether follow-up of chronic heart failure (CHF) patients, by means of a telemonitoring program, reduced mortality and rehospitalization rates (Dendale *et al.*, 2011). Daily measurements of systolic and diastolic blood pressure, heart rate and weight, were remotely collected from 80 patients. These patients at hospital discharge were provided with a set of apparatuses, through which they not only made these measurements, but also remotely availed the measurements to the medical personnel. These longitudinal measurements were recorded each day for a period of about 6 months. However, the durations of measurements depend on whether the patient get hospitalized or dropped from the study due to death or other reasons within the period which resulted in unbalanced measurements. Moreover, all biomarkers have missing values due to a variety of reasons. In addition to the biomarkers, the following patient characteristics were also collected at baseline: sex, age, heart rhythm, NTproBNP (a measure of cardiac muscle fiber stretch,the lower the value the better), patient fitness indicator(NYHA). NYHA stands for New York Heart Association functional classification of heart failure and is an indicator for patients fitness at a given moment. The scores for NYHA in the data set ranges from 1-4.The highest being the worst score. These scores are classified in to four classes based on Chul-Ho *et al.*, (2012) classifications and given in Table 1. The left ventricular ejection fraction (LVEF), is a measure of heart performance which indicates the fraction of of blood being pumped out of the ventricle within each contraction. It is usually categorized as high (or preserved_ejection) and low (or reduced_ejection) based on percentage of the blood, where a higher figure is better (Njagi *et al.*, 2013). And heart rhythm which could be classified as normal and abnormal is also available in the data set. The description of the variables are presented in Table 1 below.

Table 1: *Telemonitoring CHF Data. Variable Description.*

| Variable | category | Description |
|---|---|---|
| Predictors | | |
| Sex | 0 | Female |
| | 1 | Male |
| NYHA | 1<=NYHA<2 | classI |
| | 2<=NYHA<3 | classII |
| | 3<=NYHA<4 | classIII |
| | NYHA>= 4 | classIV |
| LVEF_status | 1(reduced_ejetion) | $\leq 45\%$ of ejection fraction |
| | 0(preserved_ejection) | $> 45\%$ of ejection fraction |
| Heartrythm | 0 | Normal heartrym |
| | 1 | Abnormal heartrym |
| Age | Age of a pateint at discharge | |
| NTproBNP | Measures of cardiac muscle fiber stretch at discharge | |
| Responses | | |
| DBP | Diastolic blood pressure (in mmHg unit) | |
| SBP | Systolic blood pressure (in mmHg unit) | |
| Weight | Weight of a patient measured (in kg times 10) | |
| Heart rate | Heart rate of a patient measured (in b.p.m) | |

# 3 Statistical Methodology

## 3.1 Generalized Additive Models

Although linear models have been the dominant approach for the analysis of longitudinal data when the outcome is continuous, in many applications the pattern of change is more faithfully characterized by a function that is non-linear in the parameters (Fitzmaurice *et al.*, 2009). In other settings, parametric models for longitudinal data are not sufficiently flexible to adequately capture the complex patterns of change in the outcome and their relationships to covariates. In such circumstances, it is fundamental to use flexible techniques such as smoothing splines (Ruppert *et al.*, 2003). Moreover, parametric regression models all assume a linear (or some parametric) form for the covariate effects and such assumption is too restrictive for many practical applications. This restrictions led to the development of nonparametric and semi-parametric regression methods, within which the linear or parametric form of (some of) the covariates are replaced by a flexible function (Fitzmaurice *et.al*, 2009). Generalized additive models (GAMs) are a nonparametric extension of GLMs, used often for the case when you have no a priori reason for choosing a particular response function (such as linear, quadratic, etc.) and want the data to show us an appropriate functional form rather than imposing some rigid parametric assumption (Wood, 2006). It extends generalized linear models (GLMs) by replacing the linear predictor with an additive predictor composed of a sum of smooth functions (Hastie and Tibshirani, 1986, 1990).

The GAM models provide a great deal of flexiblity in describing the response-predictor relationships (Wood, 2006). However, the flexibility and convenience of the model comes at the cost of two new theoretical problems. It is necessary both to represent the smooth functions in some way and to choose how smooth they should be. It requires solving 3 problems not encountered in linear modelling:(1) the smooth function has to be represented some how, (2) the degree of smoothness of the function must be controllable and (3) the amount of smoothness most appropriate should be selectable in a data-driven way (Wood, 2006). In the coming section, we will gently introduce GAM formulation and its estimation procedures as well as the choice of bases.

## 3.2 GAM formulation and base selection:

The general structure of generalized additive model is given by:

$$g(E(y_i)) = X_i^* \boldsymbol{\theta} + \sum_j f_j(x_j) \tag{3.1}$$

where $Y_i$ is the response variable, assumed to belong to some exponential famiy. g is a smooth monotonic link function, $X_i^*$ is the $i^{th}$ row of the model matrix for any strictly parametric model components, and $\boldsymbol{\theta}$ is the corresponding parameter vector. The $f_j$ are smooth functions of covariates $x_j$, which may be vector covariates. The $f_j$ are subject to identifiability constraints, typically that $\sum_j f_j(x_j)=0 \ \forall j$. The model can further be extended by including extra random effect terms to arrive at the generalized additive mixed model (GAMM)( Lin and Zhang, 1999). The first step in GAM estimation is to represent the smooth terms in model (3.1) using spline bases with associated penalties (Marx and Eilers, 1998;Wood, 2006). Each smooth term is represented as

$$f_j(x_j) = \sum_{k=1}^{K_j} \beta_{jk} b_{jk}(x_j)$$

where the $b_{jk}(xj)$ are known basis functions, chosen to have convenient properties, while the $\beta_{jk}$ are unknown coefficients, to be estimated. Given bases for each smooth term, model (3.1 ), can be re-written as a GLM, $g(E(y_i) = \boldsymbol{X_i}\boldsymbol{\beta}$, where $\boldsymbol{X}$ includes the columns of $X^*$ and columns representing the basis functions evaluated at the covariate values, while $\beta$ contains $\boldsymbol{\theta^*}$ and all the smooth coefficient vectors, $\beta_j$. The detail procedures for setting up GAM as GLM is found in Wood (2006). Unlike GLM, GAM are usually estimated by penalized likelihood maximization, where penalties are designed to suppress overly wiggly estimates of $f_j$ terms. This is the idea behind the penalized regression approach of GAM estimation (Wood, 2006). However, before GAM estimation, for each smooth function in the model, a basis has to be chosen. A basis can be seen as a way of defining the space of functions of which $f$ (or a close approximation to it) is an element (Wood, 2006). There are several spline bases to choose from. The next section introduces some of these bases and the estimation procedures for the choosen splines base for our analysis.

6

## 3.3  Bases and knots for smoothing splines

Spline functions are piecewise polynomials, with the polynomial pieces joining at the knots and fulfilling continuity conditions for the spline itself and some of its derivatives (Costa, 2008). In practice both the types of base and the number knots and its location need care (Ruppert $et,al.$, 2003). Several alternatives exist for the choice of the basis functions $b_{ji}$, such as; cubic regression splines(CRS), cyclic cubic regression splines, truncated power bases (TP), P-splines, p-splines with shrinkage, cyclic p-splines, B-splines, thin plate splines and thin plate regression splines(TPRS) to mention some. Here we will introduce some of them. The detail descriptions along with the merits and demerits of the these bases are found in Wood (2006).

**Cubic spline bases**: the most commonly used smoothing spline is the natural cubic smoothing spline. The natural cubic spline arises as the solution of the penalized residual sum of squares criterion (Hastie $et\ al.$, 2009). It interpolates the data, and yields a linear fit if the smoothing parameter $\lambda \to \infty$. Because for smoothing splines the number of smoothing parameters to be estimated is as large as the number of unique observations, they are computationally intensive (wood, 2006). On the other hand, for data sets with a large number of distinct measurement times, or with several spline terms in the model, the cubic spline mixed model generates a large number of random spline effects with a dense design matrix $Z_s$. Hence, solution of the mixed-model equations may then require a large amount of computer workspace and processing time (Ruppert $et\ al.$, 2003).

**The truncated power and B-Splines**: Truncated power bases are useful for understanding the mechanics of spline-based regression, and they can be used in practice if the knots are selected carefully or a penalized fit is used. However, the truncated power bases have the practical disadvantage that they are far from orthogonal. This can sometimes lead to numerical instability when there is a large number of knots and the penalty parameter $\lambda$ is small (or zero in the case of ordinary least squares). Therefore, in practice, especially for OLS fitting, it is advisable to work with equivalent bases with more stable numerical properties (Ruppert $et\ al.$, 2003). The most common choice is the B-spline basis. The B-spline basis is popular and widely used because of its sparse, local form and good computational properties in terms of matrix

7

inversion ( Fitzmaurice *et al.*, 2009). B-splines are also useful in constructing P-splines (Wood, 2006).

**P-splines**: P- splines are low rank smoothers which were proposed by Eilers and Marx (1996). They are low range splines in which the number of knots are much lower than the dimension of the data. They also relax the importance of the localization and the number of knots. Moreover, P-splines are extremely easy to set up and use, and allow a good deal of flexibility in that any order of penalty(measures of wiggleness) can be combined with any order of B-spline basis, as the user sees fit (Ruppert *et al.*, 2003). However, the simplicity is somewhat diminished if uneven knot spacing is required, and that the penalties are less easy to interpret in terms of the properties of the fitted smooth, than the more usual spline penalties (Wood, 2006). The focus in this thesis on p-splines using B-splines basis. In the next section we will give representation of B-spline and P-splines. With respect of the number of knots in most of the situations, the suggestion is to use a moderately large number of equally-spaced knots. The first goal for any algorithm in selecting number of knots $(K)$ is to make certain that $K$ is sufficiently large to fit the data. The second goal is to choose $K$ not so large that computation time is excessive (Durban *et al.*, 2005).

### 3.3.1   Construction of P-splines

P-splines are low rank smoothers using a B-spline basis, usually defined on evenly spaced knots, and a difference penalty applied directly to the parameters, $\beta_i$, to control function wiggliness (Wood, 2006). The B-spline basis is appealing because the basis functions are strictly local, each basis function is only non-zero over the intervals between m + 3 adjacent knots, where m + 1 is the order of the basis (Wood, 2006). To define a k parameter B-spline basis, we need to define $k + m + 1$ knots, $x_1 < x_2 < ... < x_{k+m+1}$, where the interval over which the spline is to be evaluated lies within $[x_{m+2}, x_k]$ (so that the first and last $m+1$ knot locations are essentially arbitrary). Consider a flexible regression model

$$\boldsymbol{y_i} = \boldsymbol{f}(x_i) + \epsilon_i \qquad \epsilon_i \sim N(0, \sigma^2) \qquad i = 1, ..., N \qquad (3.2)$$

8

where $\mathbf{y_i}$ is the response variable for the observations $i = 1, ..., N$ and $f(.)$ is smooth function of covariate $x$.

To estimate the function $f(.)$, it is assumed that this function can be represented by a linear combination of known basis functions $B_i$ . An $(m+1)^{th}$ order spline can then be represented as

$$f(x) = \sum_{i=1}^{k} B_i^m(x)\beta_i$$

where $\boldsymbol{\beta_i} = [\beta_1....\beta_k]^T$ is a vector of unknown regression coefficients .The $B_i$ are the B-spline basis functions, which most conveniently defined recursively as follows:

$$B_i^m(x) = \frac{x - x_i}{x_{i+m+1} - x_i} B_i^{m-1}(x) + \frac{x_{i+m+2} - x}{x_{i+m+2} - x_{i+1}} B_{i+1}^{m-1}(x) \quad i = 1, ...k$$

$$B_i^{-1} = \begin{cases} 1 & x_i < x < x_{xi+1} \\ 0 & otherwise \end{cases}$$

Under this representation the model parameter $\theta_i$ can be easily estimated using ordinary least squares;

$$\hat{\beta} = \boldsymbol{B}^T \boldsymbol{B}^{-1} \boldsymbol{B}^T \boldsymbol{y}$$

where

$$\boldsymbol{B} = \begin{vmatrix} \boldsymbol{B}_1(x_1) & ... & \boldsymbol{B}_K(x_1) \\ \vdots & \ddots & \vdots \\ \boldsymbol{B}_1(x_N) & ... & \boldsymbol{B}_K(x_N) \end{vmatrix}$$

and $\mathbf{y} = (y_1....y_N)$ and the penalization is discrete and directly penalizes the coefficients instead of the whole curve which reduces the dimensional problem.

## 3.4 GAM estimation via P-IRLS

We have stated that in regression spline, the number and also the location of the knots have a great impact on the final estimates. When the number of knots increases the estimated curve (in comparison with the true curve) becomes too wiggly, meaning that the data are over fitted (Wood, 2006). To solve this problem, O'Sullivan (1986) introduced the idea of penalized splines, where a smoothness penalty is added to the least squares criterion while estimating the regression coefficients $\boldsymbol{\beta}$. For example, the penalization based on the second derivative function of $f$ of penalized splines for model (3.1) is fitted by minimizing the penalized sum of squares:

$$(\boldsymbol{y} - \boldsymbol{B}\boldsymbol{\beta})^T(\boldsymbol{y} - \boldsymbol{B}\boldsymbol{\beta}) + \lambda \int f''(x)^2 \partial x \tag{3.3}$$

where $\lambda$ is the smoothing parameter which controls the trade-off between fidelity to the data and roughness of the function estimate. The trade off between model fit and model smoothness is controlled by the smoothing parameter $\lambda$. $\lambda \to \infty$ leads to a straight line estimate for f, while $\lambda = 0$ results in an un-penalized regression spline estimate.

Because f is linear in the parameters, $\boldsymbol{\beta_i}$, the penalty can always be written as a quadratic form in $\boldsymbol{\beta}$

$$\int [f''(x)]^2 \partial x = \beta^T S \beta$$

Associated with each smooth function ,the measures of function wiggliness is given by $\boldsymbol{\beta_j^T} \tilde{\boldsymbol{S}}_j \boldsymbol{\beta_j}$, where $\tilde{\boldsymbol{S}}_j$ is a matrix of known coefficients. As stated in section 2.2, given bases for each smooth term, the GAM model can be re-written as a GLM. However, the generalization from GLMs to GAMs requires the development of theory for penalized regression, in order to avoid problems of overfitting (Wood, 2006). The fit of such GLM is most conveniently measured using the deviance:

$$D(\boldsymbol{\beta}) = 2\{l_{max} - l(\boldsymbol{\beta})\}\phi$$

where $l$ is the log-likelihood of the model, and $l_{max}$ is the maximum possible value for $l$ given the observed data, which is obtained by considering the MLE of a model with one parameter

10

per datum (under which the model predicted $E(y_i)$ is simply $y_i$)(Wood, 2008). $\phi$ is is a scale parameter, and the definition of D means that it can be calculated without knowledge of $\phi$. Maximizing the likelihood is equivalent to minimizing the deviance, and in several ways the deviance behaves rather like the residual sum of squares in linear modeling (Wood, 2008). If the bases used for the smooth functions, $f_j$, are large enough to be reasonably sure of avoiding misspecification, then the model will almost certainly overfit if it is estimated by minimizing the deviance (Wood, 2008 ). For this reason GAMs are estimated by minimizing

$$D(\beta) = \sum_j \lambda_j \beta^T S_j \beta$$

where $\lambda_j$ are smoothing parameters and the $S_j$ are the $\tilde{S}_j$ suitably padded with zeroes so that $\beta^T S_j \beta = \beta_j^T \tilde{S}_j \beta_j$ and $S = \sum_j \lambda_j S_j$. The $\lambda_j$ controls the smoothness of the component smooth functions. Hence, smoothness selection is about choosing values for the $\lambda_j$. Given $\lambda$ ,the penalized deviance can be minimized by penalized iteratively re-weighted least squares (P-IRLS)(Wood, 2006). However, $\lambda_j$ have to be estimated as well. The fact that $\lambda_j$ is unknown, arise a need for an iterative procedure to estimate $\beta$.

Let $V(\mu)$ be the function such that $var(y_i) = V(\mu_i)\phi$. Let $w_i$ denote any prior weights on particular data points (used to weight the component of deviance attributable to each datum). Then, $\hat{\beta}$ can be estimated by iterating the following steps to convergence.

1. Given the current $\mu^k$, calculate the pseudo data $Z^k$, and weights,$w_i{}^k$,where
   $w_i{}^k = \frac{1}{V(\mu_i{}^k)g^T(\mu_i{}^k)^2}$ and $Z_i = g^T(\mu_i{}^k)(y_i - \mu_i{}^k) + \mathbf{X}\hat{\beta}^{\mathbf{k}}$, g is the link function, $Z^k$ is a vector of pseudo-data and $W^k$ is daiagonal matrix with diagonal elements $w_i^k$.

2. Minimize
   $\|\sqrt{\mathbf{W^{[k]}}}(z^K - \mathbf{X}\beta)\|^2 + \beta^T S_j \beta$ with respect to $\beta$ to find $\hat{\beta}^{[k+1]}$.
   Evaluate the linear predictor, $\eta^{[k+1]} = X\hat{\beta}^{[k+1]}$ and the fitted values $\mu_i^{[k+1]} = g^{-1}(\eta_i^{[k+1]})$.

Hence, the penalized least squares estimator of $\beta$ is given by;

$$\hat{\beta} = X(X^T X + \lambda S)^{-1} X^T \qquad\qquad (3.4).$$

Where the influence, or the hat matrix for the model can be written $A = X(X^T W X + S)^{-1} X^T W$.

## 3.5  Degrees of Freedom and Scale parameter estimation

The effective degrees of freedom of a GAM, defined as tr(A), where A is the influence matrix described above, indicate the flexibility of the fitted model. For instance, using large values for the smoothing parameters would result in a very inflexible model (a model with very few degrees of freedom). The application of penalties reduces the model degrees of freedom. It is possible to break down the effective degrees of freedom of the model to each smooth function in the model or even to each $\hat{\beta}_i$ separately. It can be shown that the effective degrees of freedom for the model parameters in the general weighted case are given by the leading diagonal of

$$F = (X^T W X + S)^{-1} X^T W X, \tag{3.5}$$

where $S = \sum_j \lambda_j S_j$. F can also be shown to be the matrix that maps the un-penalized estimates to the penalized ones and $F_{ii}$ to measure the effective degrees of freedom of the $i^{th}$ penalized parameters.

**Residual Variance or scale parameter estimation**: For the additive model case,the residual variance($\sigma^2$) is usually estimated in a manner analogous to the linear regression case as

$$\hat{\sigma}^2 = \frac{\|y - Ay\|^2}{n - tr(\mathbf{A})} \tag{3.6}$$

whilst the scale parameter in the case of a GAM is estimated by the Pearson-like estimator as

$$\hat{\phi} = \frac{\sum_i V(\hat{\mu}_i^{-1}(y_i - \hat{\mu}_i)^2}{n - tr(\mathbf{A})} \tag{3.7}$$

## 3.6   Smoothing parameter selection

The critical point of the penalized spline smoothing is the choice of the smoothing parameter $\lambda$. In the previous section, it was stated that if a large smoothing parameter is chosen the resulting curve is very smooth, but if we choose a small smoothing parameter, the resulting estimate becomes too wiggly. Penalized likelihood maximization can only estimate model coefficients,$\beta$, given smoothing parameters $\lambda$, so it is of interest to discuss how to estimate the smoothing parameter and then this section covers some practical methods of $\lambda$ estimation.

There are several methods to choose the optimal value of of $\lambda$ such as Ordinary cross-validation (OCV), Generalized Cross Validation (Hastie and Tibshirani, 1990; Wood, 2006) or Akaike Information Criterion (Wood, 2008). The methods of Ordinary cross-validation (OCV) and GCV were introduced as automatic methods of smoothing parameter choice. OCV aimed to minimize the mean squared error of prediction (MSEP) across a set of data points. By omitting observation $y_i$, while fitting the model, using the model to predict $E(y_i)$, and repeating the procedure to all data in turn, the following estimate of OCV in the additive model case is obtained:
$$V_o = \tfrac{1}{n}\sum_i^n (y_i - \hat{\mu}_i^{[-1]})^2$$

where $\hat{\mu}_i^{[-1]}$ denotes the prediction of $E(y_i)$ obtianed by omitting $y_i$. It can be shown that estimation of $V_o$ can be estimated once using the following expression which requires fitting the original model once.

$$V_o = \frac{\sum_{i=1}^{n}(y_i-\hat{\mu}_i)^2}{n(1-A_{ii})^2}$$

However, the OCV method is computationally expensive and suffer from lack of invariance (Wood, 2006). Another disadvantage of this method was that data values with high leverage could have a large influence on the choice of $\lambda$. On the other hand, Genaralized cross validation( GCV) (Craven and Wahba, 1979) was developed to downweight these values in the estimation of the smoothing parameter. The GCV method obtains an estimate of the value $\lambda$ that minimizes the MSEP. In additive model, GCV is given by

$$V_g = \frac{n\|(y-\mu_i)\|^2}{[n-tr(\mathbf{A})]^2}$$

13

Generaization of GCV criterian into GAM can be obtained first by writting the GAM fitting objective in terms of model deviance

$$D(\boldsymbol{\beta}) + \sum_{j=1}^{m} \lambda_j \beta^T D_J \beta,$$

then the GCV score is defined as

$$V_g = \frac{nD(\hat{\beta})}{[n-tr(\mathbf{A})]^2}$$

GCV has computational advantages over OCV, and it also has advantages in terms of invariance (Wahba, 1990). Efficient algorithms are available to calculate both the OCV and the GCV estimates. However, both OCV and GCV estimates may work poorly in the presence of correlated errors, although they are thought to be reasonably robust against misspecification of the error distribution( Fitzmaurice *et al.*, 2009).

# 4 Generalized Additive Mixed Models

When repeated measurements are taken on each subject, the subject to-subject variation introduces a new source of randomness and an extension to GAMs may be necessary. Analogous to the extension of GLMs to generalized linear mixed models (GLMMs), GAMs have likewise been extended to GAMMs (generalized additive mixed models) by inclusion of random effects. A GAMM is just a GLMM in which part of the linear predictor is specified in terms of smooth functions of covariates (Lin and Zhang, 1999). For example, an Additive Mixed Model has a structure something like

$$y_i = \boldsymbol{X_i}\boldsymbol{\beta} + f_1(x_{1i}) + f_1(x_{2i}, x_{2i}) + \cdots + \boldsymbol{Z_i}\boldsymbol{b} + \epsilon_i \tag{4.1}$$

where $y_i$ is a univariate response; $\beta$ is a vector of fixed parameters; $\mathbf{X}_i$ is a row of a fixed effects model matrix; the $f_j$s' are smooth functions of covariates $x_k$; $\boldsymbol{Z_i}$ is a row of a random effects model matrix; $\boldsymbol{b} \sim \boldsymbol{N(0, \psi)}$ is a vector of random effects coefficients, with unknown positive definite covariance matrix $\boldsymbol{\epsilon_i} \sim \boldsymbol{N(0, \Sigma)}$ is a residual error vector, with $i^{th}$ element $\boldsymbol{\epsilon_i}$, and covariance matrix $\Sigma$ which is usually assumed to have some simple pattern.

In section 3.4, it was said that moving from GLMs to GAMs requires the development of theory for penalized regression, in order to avoid overfitting, but GLMM methods require no adjustment in order to cope with GAMMs: it is possible to write any of the penalized regression smoothers as components of a mixed model, while treating their smoothing parameters as variance component parameters, to be estimated by Likelihood, REML or PQL methods (Wood, 2006). The detail procedures of estimation and the inferencial paradigm with respect GAMM is found in Wood (2006).

## 4.1 Model Selection

In GAM frame work, model can be compared using the GCV deviance or AIC (Wood, 2006). For each model, GCV or AIC would be minimized over $\lambda$ to obtain the model GCV or AIC. Then models with small GCV or AIC would be considered best. All models with reasonably small GCV or AIC values should be considered as potentially appropriate and evaluated according to their simplicity and scientific relevance. The idea behind the AIC is to penalize the loglikelihood

15

with the number of parameters, namely AIC=-2LL+2p with p the number of parameters in the model (fixed effects and variance components). The usage of the AIC in this form may not be appropriate in case of semi parametric models (Maringwa *et al.*, 2008c) and instead an adjusted AIC, should be used. The penalty term of the adjusted AIC takes the effective number of parameters into account, which generally is higher than p because the smoothing is accounted for. In Sectoin (3.6) we have seen how to determine the effective degrees of freedom(or effective number of parameter) by using the design matrix of fixed and random effects corresponding to the penalized spline. Thus the effective number of parameters$(E_p)$ for the generalized model is (Ruppert *et al.*, 2003) given by equation (3.5), then adjusted AIC is given by $AICadj = -2LL + 2Ep$.

## 4.2    Hyphothesis Testing

Next to model selection, formal tests are required. Hastie and Tibshirani (1990) describe approximate F-tests for generalized additive models based on deviances. One may also wish to test the hypothesis that the simpler of two nested models is correct versus that the larger is correct. A comparison of the likelihoods evaluated at the penalized maximum likelihood estimates might be given, by the following test statistic:

$$\lambda = 2\|l(\hat{\eta}_l) - l(\hat{\eta}_0)\|$$

where $\hat{\eta}_0$ and $\hat{\eta}_l$ are the estimated linear predictors for the null and alternative models respectively. Treating the smoothing parameters as known rather than estimated, the test statistics is roughly approximated that under the null, $\lambda \sim \chi^2_{EDF_l - EDF_0}$, where $EDF_1$ and $EDF_0$ representing the respective effective degrees of freedom.

To test for the need of covariance matrix parameters, likelihood ratio test based on mixture of chisquares $(\frac{1}{2}\chi^2_s + \frac{1}{2}\chi^2_{s+1})$ can be used. It is an asymptotic distribution where $s$ is the number of fixed effects parameters constrained under the null hypothesis.

For inference regarding fixed effects, the bayesian credible intervals was used. From the fact that we rely on some prior beliefs about parameter $\lambda$, the distribution for GAM results in bayesian

frame works (Wood, 2006). Hence, for constructions of confidence bands the bayesian uncertainity approach was used. The detail descriptions of the distributional results and derivation of confidence bands for GAMs and also for GAMMs are found in Wood (2006).

# 5 Software

In our report, SAS 9.3 was used to get summary statistics and plots. The mgcv package within R 3.2 was used for all inferential analyses. Specifically, the gamm function within the mgcv package was used for GAMM fitting. The function implements GAMMs in two phases: first calling lme or glmmPQL to estimate the model in the AMM and GAMM case, respectively, and then transforming the returned object into a gam object, so that GAM-related inferences can be drawn. The estimated plots along with their credible intervals are obtained from gam component of the model, while inference for model diagnostics are obtained from lme component of the models. Selected codes are shown in Appendix: E for both SAS and R. All statistical tests were conducted at 5% level of significance.

# 6  Results

In this section, the data introduced in Section 1.3 are analyzed and results of the analysis based on GAM and its extension with application of the penalized spline (P-splines) will be discussed. Recall that the aim of this study is to come up with a flexible technique to model postdischarge profiles of telemonitered chronic heart failure patients; i.e investigating the functional relationships of the longitudinal biomarkers measurements and time, and to see if there is a difference in the evolutions of the profiles across the patient's baseline characteristics described in Table 1. Analysis of each outcome will be presented separately. First, diastolic blood pressure by taking into account the possible covariates will be discussed. Then, systolic blood pressure, weight and heart rate measurements will be analyzed in the respective order.

## 6.1  Exploratory Analysis

Prior to analysis, the data were examined for possible problems. A large number of missing observations was observed in the data in general and in each outcome in particular for unknown reason of missingness. As mentioned in section 1.3 the study involved 80 subjects whose biomarkers were measured repeatedly during 6 months of follow-up after hospitalization. The number of biomarker observations per subject varied from 0 to 186, resulting in a total of 9767, 9740, 10183 and 9744 observations for diastolic, systolic, weight and heart rate respectively. Table 2 shows the summary statistics of the four outcomes. Out of 80 patients, one had no measurements for DBP and SBP while 4 had no measurements for measures of cardiac muscle fiber stretch (NTproBNP). The value of DBP ranges from 31 to 125 with mean of 71.4, while SBP ranges from 74 to 205 with mean of 123.83 in mmHg units. The ranges observed in weight measurements was considerable, ranging from 350.5 to 1478 with mean of 766 in units of 10xkgs. The measurements for heart rate ranges from 40 to 134 with mean 69.8 and variance of 164.4 in beats per minute ( b.p.m). In the data there were 30 female and 50 male patients. With respect to the left ventricular ejection fraction(LVEF), 23 patients had high ejection fraction ($> 45\%$) while 57 patients had the low ($< 45\%$) figure. By heart rhythm, 44 patients had normal while

36 classified as abnormal. The distribution of patients by New York Heart Association classes (NYHA) is highly unbalanced. 34 pateints were categorized to class III followed by class IV with 27 patients. Class I and class II had 5, 14 patients respectively. The median age of patients was 77 whith minimum 46 and maximum age of 95.

Table 2: *Summary statistics*

| Variables | N | NMISS | Minimum | Mean | Maximum | Variance |
|---:|---|---|---:|---|---|---|
| DBP | 9767 | 2566 | 31 | 71.35 | 125 | 146.61 |
| SBP | 9740 | 2593 | 74 | 123.83 | 205 | 385.47 |
| Weight | 10183 | 2150 | 390 | 766.06 | 1478 | 28007.49 |
| Heart rate | 9744 | 2589 | 40 | 69.73 | 134 | 164.41 |

### 6.1.1 Diastolic Blood Pressure

**Individual Profiles**

Figure 1 shows the individual diastolic postdischarge profiles for randomly selected patients. The profiles appeared indistinguishable with up and down curvatures. It demonestrates cyclic trends which makes the assumptions of linear trend questionable. Moreover, it shows high within and between variability which suggest the plausibility of random effects for the model to be considerd.
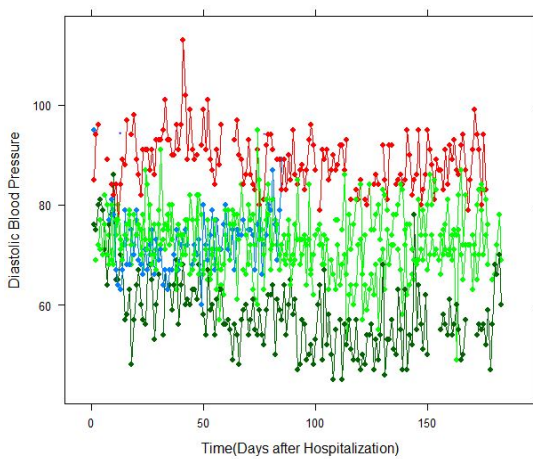


Figure 1: *Randomly selected individual profiles of daistolic blood pressure.*

Figure 2: *Average evolution of diastolic blood pressure.*

In order to be able to further visualize the overall evolution of diastolic blood pressure, the overall mean evolution for the diastolic blood pressure was considered and can be viewed in Figure 2. The curve demonestrates up and down trend through out the period. The mean profiles by level of covariates are also given in Figure 3 (a-e). The plot shows distinguished evolution of the average diastolic blood pressure accross the categories of age, but the difference gradually vanish. Similarly, there is a difference in the evolution across class of NYHA, where the trends for patients in the first three classes evolved differently from those in the fourth class. For Sex, LVEF and Heart rhythm, the profiles are indistinguishable except at some points in time. The observed individual and mean profiles were not so clear for one to describe the post discharge profiles of diastolic blood pressure. Hence, the profiles need to be smoothed in appropriate ways to describe the longitudinal trends and there by differentiate the evolutions among groups of patients in the study.

(a)



(b)



(c)



(d)



20

(e)



Figure 3: *Average evolution of diastolic blood pressure by levels of covariates. (a) gender, (b) LVEF, (c) heart rhythm,(d) age category and (e) NYHA.*

### 6.1.2 Systolic Blood Pressure

The subject specific longitudinal systolic profiles for randomly selected patients are presented in Figure 4, where the within-patient variability appeared substantial. In Figure 5 the overall average evolution is shown. The trend is increasing although it shares the wiggly curvature observed in DBP profiles of Figure 2. The systolic mean profiles by level of covariates are displayed in Figure 23 (in appendix B). The plot shows distingiushed evolution of postdischarge systolic blood pressure profiles for males and females as well as across LVEF levels. The average evolution of females remained higher than males during the period. Similarly, the profile for patients who had high ejection fraction (LVEF > 45 %) is clearly distinguished from those who had low ejection fraction. For the NYHA, the evolution is some what differ accross classes where patients in the fourth and second classes had a lower evolution range compared to the first and the third classes. Specifically, for patients in the fourth class the trend remained below the pulled average systolic value for most of the time. For heart rhythm and age category, there was no clear difference in the trends except at some intervals. All in all, the observed individual and average systolic profiles implied that a suitable model would have to be adopted to capture the

21

oscillations in a smooth rather than constant fashion.



Figure 4: *Randomly selected individual profiles of systolic blood pressure.*



Figure 5: *Average evolution of systolic blood pressure.*
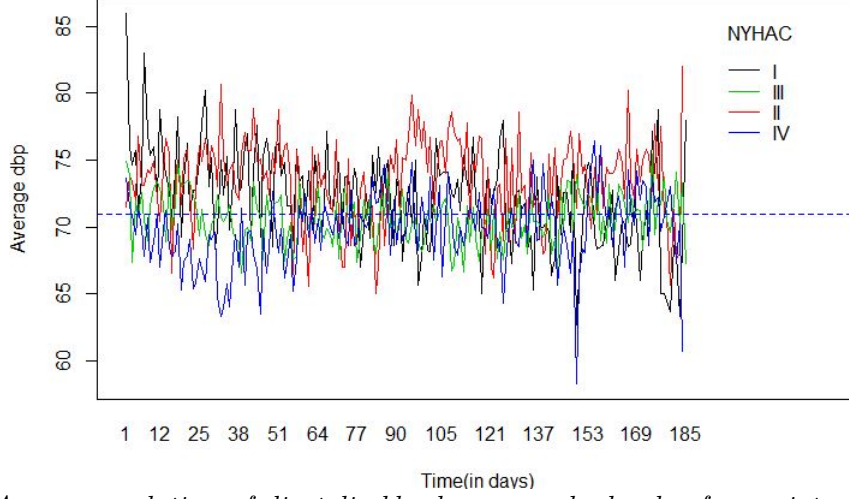
### 6.1.3 Weight

The individual and average weight trends are presented in Figure 6 and Figure 7 respectively. The individual profiles appeared in a different range of evolution suggesting there was a high between variability but more or less constant within fluctuation unlike systolic and diastolic profiles. The overall mean profile fluctuates around the pooled average for most of the durations, but it sharply decreased at the end. Furthermore, the mean profiles by levels of covariates are displayed in Figure 24 (in appendix B). It shows distinguished evolution of weight measurements conditioned on the categories of covariates. For patients with worst score of NYHA (class IV), around the first month of post discharge, their average profile remained fluctuating above the pooled average value, but dropped immediately and start evolving below the average value for the rest of the duration. Once again Figure 24 can give a clue to consider all covariates in studying postdischarge profiles of weight and the need of smoothed profiles.

Figure 6: *Individual profiles of Weight.*



Figure 7: *Average evolution of Weight.*

### 6.1.4 Heart rate

In Figure 8 and Figure 9 individual and average profiles for heart rate measurements are presented respectively. The within and between Variability is evident from the individual and mean postdischarge profiles, with mean profile decreasing over time . On the other hand, Figure 25 ( in appendix B) shows mean heart rate profiles by covariates. It can be seen that there was a difference in evolution across levels of covariates. In general the observed curvatures suggest the need for a flexible model, which would be able to capture the functional dependence of heart rate on time. As the aforementioned outcome variables, we need to study the profiles in a better way such that one can easily describe the trend.



Figure 8: *Randomly selected individual profiles of heart rate measurements*



Figure 9: *Average evolution of heart rate*

23

# 7 Modelling The Trends

## 7.1 GAMM Application to Chronic heart failure data

In this section, we apply the GAMM model (4.1) to analyze the longitudinal biomarkers. The interest is in estimating the time profile of these biomarker measurements as well as investigating the effects of covariate, on the trends. All variables, except age and NTproBNP are dichotomous (LVEF (high=1, low=0), gender (male=1, female=0), heart rhythm (normal=0,abnormal=1), and NYHA_classes (I, II, III, IV). For each of categorical variables three different models were considered. Model 1 which assumes the overall trend of the biomarkers (or common effect of time with out adjusting for covariates). In model 2, the trends for the groups differ only by constant, and model 3 which allows different trends for the groups (the difference curve). Later on, different modifications were made on model 3 by varying the structures of random effects.

The general formulations of the models are given as :
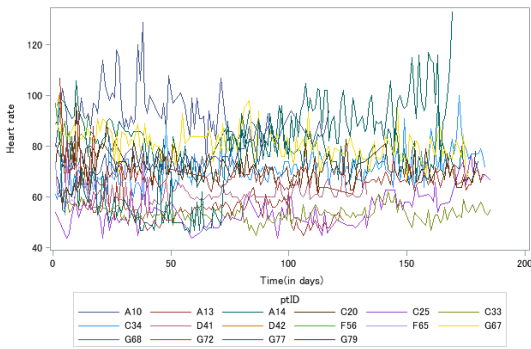
$$Y_{ij} = f(t_j) + b_i + \epsilon_{ij} \tag{1}$$

$$Y_{ij} = \beta_0 + \beta_1 X_k + f_1(t_j) + b_i + \epsilon_{ij} \tag{2}$$

$$Y_{ij} = f_1(t_j) + f_2(t_j) X_k + b_i + \epsilon_{ij} \tag{3}$$

where $Y_{ij}$ is the response of patient $i$ measured at occasion $j$ ,$i=1, \cdots , N$, $j=1, \cdots , n_i$, $f(.)$ is a smooth function which reflects response trend, $t_j$ is the time of interest in days, $b_i$ is random intercept, $\epsilon_{ij}$ is the random error of the $j^{th}$ response measurement for $i^{th}$ patient, $X_k$ is the $k^{th}$ level of the covariate representing Sex, LVEF, heart rhythm and NYHA in their respective analysis. $\beta_0$ is the intercept for the reference group and $\beta_1$ represents the fixed effect parameter for $k^{th}$ level of the covariate .

$$b_i \sim N(0, D)$$

$$\epsilon_{ij} \sim N(0, \sigma^2)$$

The continuous covariates, age and NTproBNP were also included in the model, but in parametric form.

### 7.1.1 Comparison of spline bases

Although we proposed to use p-splines, it was of interest to compare the performance of some of the bases. According to (Ruppert *et al.*, 2003), in principle, a change of basis does not change the fit though some bases are more numerically stable and allow computation of a fit with greater accuracy. Besides numerical stability, reasons for selecting one basis over another are ease of implementation (especially of penalties) and interpretability. The latter consideration is usually not too important since one is generally interested only in the fit, not the estimated coefficients. Prior to fitting tha actual model, we carried out comparison of some of spline bases. As stated in section 3.2 there are several spline bases such as truncated polynomial basis, cubic regression splines (CRS), cyclic cubic regression spline, thin plate regression spline (TPRS), p-splines, and shrinkage smoothers to mention some. We have also discussed the theoretical properties of some of these bases. Now we will explore a number of these options in terms of the smoothing ability of the model, and its fit. For illustration purpose, we considered GAMM for the diastolic blood pressure with only the smoothed effect of time i.e model (1).

First, model (1) was fitted using cubic regression splines (CRS), p-splines (ps), thin plate regression splines (TPRS), cubic regression splines with shrinkage (cs), thin plate regression splines with shrinkage (ts), cyclic cubic regression splines (cc), cyclic p-plines (cp) and the tensor product smooth with p-splines bases. In Figure 10, the result of model 1 for the smoothed effect of time on diastolic blood pressure is presented. The y-axis represents the estimated effective degrees of freedom (edf) while the x-axis represents the time in days. It can be seen that the trend for DBP seems more or less the same under each option of bases. The different smoothing bases produce almost indistinguishable estimated curves. The estimated effective degrees of freedom (edf) range from 3.5 to 5, where the lower edf indicates higher penality. The solid line is the estimated effect where the dashed lines are the confidence bands. The horizontal dashed line shows the region where the effect is zero.

In addition to the use of plots, the bases were formally compared using AIC . Table 3 displays the AIC and loglikelihood values of the fitted models. The model with cubic regression splines

Figure 10: *Graphical comparison of smoothing bases for the random intercept model with smoothed time effect fitted to diastolic blood pressure.*

and p-splines gave the smallest AIC which is equal to 66969.41 and 66971.14 respectively.

By adjusting AIC for the respective total estimated effective degrees of freedom (5.720 and 5.230) using $AICadj = -2LL + 2Ep.$, for p and cubic regression splines the adjusted AIC became 66970.84 and 66971.54 respectively. This suggested that the two bases do the job equally well.

The above considerd random-intercept model only assumes a shift in subject-specific profiles, which is a restrictive assumption. So it is of interest to comapre bases by considering a more complex models, for example, including subject-specific intercepts and slopes. Hence, to see the long term effect while accounting subject to subject variation, comparison was done by extending model (1), where random slope is added to the model . As shown in Table 3 the inclusion of random slope improved the model loglikelihood and resulted in smallest AIC compared to the random intercept model. Moreover, the AIC for the two bases, cubic regression and p-splines were almost equal with 65820.18 and 65820.29 respectively.

26

Table 3: *Comparison of bases using Loglikelihood and Akaike Information Criterion. RI is for random intercepts model while RS is the model with random intercept and slopes.*

| Model | | CRS | PS | TPRS | cc | cp | cs | ts | te with ps |
|---|---|---|---|---|---|---|---|---|---|
| RI | AIC | **66969.41** | **66971.14** | 66973.64 | 66984.2 | 66984.49 | 66974.28 | 66976.83 | 66972.99 |
| | -loglik | -33479.71 | -33480.57 | -33481.82 | -33488.10 | -33488.24 | -33483.14 | -33484.42 | -33480.99 |
| RS | AIC | **65820.18** | **65820.29** | 65821.75 | 65821.34 | 65822.54 | 65822.18 | 65841.44 | 65823.69 |
| | -loglik | -32902.09 | -32902.24 | -32902.87 | -32903.67 | -32903.77 | -32904.09 | -32913.72 | -32904.84 |

Lastly, we compared the fit performance of some smoothing bases by using model (1). But, now with out the random effects term , and where the mean of diastolic and systolic blood pressure were fitted as a function of smoothed effect of time (in days). The results of the models are presented in Figure 11 and 12. The figure shows the observed mean of DBP and SBP overlaid by the estimated values under the different bases. The estimated trend lines from cubic regression (green solid line) and p-splines (red solid line) bases approximated the observed means well. In conclusion, the comparison revealed that for the data at hand, both cubic regression and p-splines provided a good deal of flexibility to the profiles. However, because of the computational issue and for other reasons stated in section (3.3), p-spline was adopted for further analysis.

In the coming section we will illustrate the GAMM model for each biomarker by considering the influence of baseline patient characteristics one at a time using p-splines. But, before applying the mixed model structures, the simple GAM models (with out random effects) were fitted for the four outcomes as a smoothed functions of only time in the models. The results are illustrated in Figure 26 in appendix B. It can be seen that the estimated mean curve for diastolic showed a fluctuating trend. The systolic curve is increasing while it resulted in a decreasing trajectory for heart rate. For weight, it initailly declined then followed by constant evolution and at the end it dropped quickly.

Figure 11: *Average DBP profiles overlaid by smoothed fitted averages.*



Figure 12: *Average SBP profiles overlaid by smoothed fitted averages.*

## 7.2 Diastolic Blood Pressure

The three models described in section 7.1 were considered by taking the influence of each covariate at a time. For example, model (3) for the influence of gender on diastolic blood pressure is given by :

$$DBPij = f_1(t_j) + f_2(t_j)Sex_i + b_i + \epsilon_{ij}$$

For $i = 1, 2, ......, N, j = 1; 2, ...., ni$, $DBP_{ij}$ is the diastolic measurement of $i^{th}$ pateint at $j^{th}$ day. where $f_1$ is a smooth function representing the trajectory for male patients, $f_2$ the smooth trajectory of female patients. The aim is to see the long term effect of gender so that we have used an interaction model in which the categorical factor (gender) interacts with a continuous factor (day). Both $f_1$ and $f_2$ were represented using p- splines. The function $\boldsymbol{f(.)}$ is estimated by $\boldsymbol{f(t_j)} = \boldsymbol{B^j\theta}$ where $\boldsymbol{B^j}$ is the B-spline basis for the function $f(t_j)$ and $\boldsymbol{\theta}$ is the fixed effect parameter . In this model $bi$ represents the subject specific random intercepts.

Table 4 contains the gam components of the fitted model where edf represent the estimated effective degrees of freedom. The diastolic trajectories for both male and female patients were estimated with a significant smoothed trends. The AIC values for model (1) and model (3) were 66971.14 and 66956.18 respectively, which favored model (3) over model (1). On the other hand,

in comparison of the long term smooth effect of gender, model (3) with a model which assumes constant effect of gender, model (2) (with AIC=66972.87), clearly supported the former over model (2). The resulting plot for model (3) is also displayed in the upper two panels of Figure 13. A straight line, corresponding to 2 degrees of freedom, was estimated for the female trajectory, while the effect of male was estimated as a smooth curve with 4.253 degrees of freedom. The total degrees of freedom, the sum of these two plus one degree of freedom for the model intercept resulted in 7 degrees of freedom. The solid lines/curves on the plots represent the estimated effects while the dashed lines corresponds to the 95% confidence limits (strictly Bayesian credible intervals). The rug plots, along the bottom of each plot, show the values of the covariates of each smooth, while the number in each y-axis caption is the effective degrees of freedom of the term being plotted. From the figure we can infer that the diastolic profiles for male patients is decreasing with some fluctuations while for female patients it had an increasing trend.

Table 4: *Gam components of the model fitted to diastolic blood pressure in assessing gender effect (RI is for random intercepts and RS for random slopes).*

|  | Smoothed.term | edf | Ref.df | F | p.value |
|---|---|---|---|---|---|
| Model (1) | s(day) | 4.235 | 4.235 | 9.605 | <0.0001 |
|  | R-sq.(adj) | = 0.0001 | AIC | =66971.14 | |
| Model (2) | s(day) | 4.253 | 4.253 | 15.04 | <0.0001 |
|  | R-sq.(adj) | = 0.0013 | AIC | = 66972.10 | |
| Model(3) with RI | s(day)*(sex=M) | 4.253 | 4.253 | 15.04 | <0.0001 |
|  | s(day)*(sex=F) | 2.000 | 2.000 | 10.49 | <0.0001 |
|  | R-sq.(adj) | = 0.0028 | AIC | =66956.18 | |
| Model(3) with RS | s(day)*(sex=M) | 5.504 | 5.504 | 13.892 | <0.0001 |
|  | s(day)*(sex=F) | 6.664 | 6.664 | 7.851 | <0.0001 |
|  | R-sq.(adj) | = 0.0028 | AIC | = 65789.27 | |

Model (3) was further extended to a more elaborative structure by allowing subject specific random slopes. To test if the random slope assumption is appropriate for the model, the model was fitted with and with out random slopes using the following hypothesis of interest:

$$H_0 : \sigma^2_{U2} = 0 \quad ; \quad H_1 : \sigma^2_{U2} > 0$$

29

Where $\sigma_{U2}^2$ is the variance of random slopes. As the null distribution of the test statistics, a mixture of chi-square approximation ($\frac{1}{2}\chi_0^2 + \frac{1}{2}\chi_1^2$) was used. The $-2Restricted\ log\text{-}likelihood$ was 1172.909 and the p-value was <0.0001, highly significant, so there is evidence to reject the null hypothesis. That means the random slope would be appropriate for the model. The residual error variance from the model with random slopes was estimated as $\sigma_e^2 = 45.99$ and the variance parameter of random intercept was $\sigma_{u_1}^2 = 87.24$, which indicated a high patient to patient variability at baseline. The variance parameter of the random slope was $\sigma_{u_2}^2 = 0.0033$. Negative correlation was observed between the random intercept and slopes, which was estimated as $\gamma = -0.755$, it showed that patients who started high do not remain high and vice versa. The plot for gam components of model (3) with the random slopes is displayed in the lower two panels of Figure 15, where the female trajectory showed some smoothed curvature compared to the straight line of model without random slopes. The smoothed diastolic curve for male patients showed both decreasing and increasing trend during the period. For female patients, the curve initially increased and then started to evolve more or less in a constant fashion. In general the estimated curves under the random slope model appear to describe the mean evolution rather well. NTproBNP and age were also included in the model in parametric form however, NTproBNP was not significant at 5% significance level while age had a significant impact on the evolution. For model diagnostics, a plot of standardized residuals versus time was plotted for model 3 with random slopes structure and shown in Figure 14. Residuals were observed to scatter without showing any particular trends, which is an indication that the selected model approximated the data well.

Figure 13: *Estimated smooth terms of the model for effect of gender on DBP. The upper two panels are from model(3) with random intercept only and the lower are obtained from random slopes and intercepts model*



Figure 14: *Model diagnostic for gender effect on diastolic blood pressure resulted from the model with random slopes.*

**LVEF, Heart rhythm and NYHA**: So far we have looked at the influence of gender on diastolic longitudinal profiles. In what follows we will study the long term effect of LVEF, heart rhythm and NYHA respectively. For each predictors, three models, i.e model (1), model (2) and

model (3) were considered. From comparison of the model with constant effect of the covariates model (2) with model (3), the latter resulted in smallest AIC and prefered over the former which happened true for the three covariates (Table 9 in appendix A). Hence, for all predictors model (3) was choosen as optimal model. The numerical output for the estimated effects from model (3) is presented in Table 5. From the table, we can see that the longitudinal diastolic profile for LVEF was estimated as a smooth curve with 6 and 4.3 degrees of freedom for patients who had high and low level of ejection fraction respectively. For normal and abnormal levels of heart rhythm the trajectories were estimated non linearly with 5.84 and 6.45 degrees of freedom. Similarly for NYHA classes except for the patients i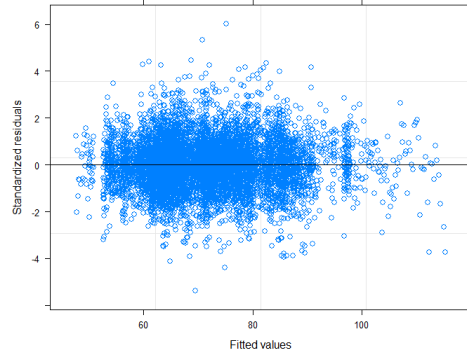n class II, the curve had a significant non linear trend. The gam components of the model for left ventricle ejection fraction and heart rhythm are plotted and presented in Figure 15. The curves for LVEF evolved differently for the two categories early in time after hospitalization. Patients with high ejection fraction showed a decreasing trend while those with low level had an increasing trend. Similarly, non linear trends are observed for patients with normal and abnoraml heart rhythm. In Figure 16 the estimated diastolic curve by NYHA are presented. Patients in the first (I) category had a sharp decreasing profiles with 1 degrees of freedom. Patients in class II had straight line trajectory with 2 degrees of freedom while the rest class showed a curvature trend; the third (III) class showed an increasing smoothed trend. Patients in the fourth (IV) class showed a different evolution phases, where the initial decreasing phase followed by along increasing period before reaching the plateau at about day 100 and then showed some fluctuation. The random slope model was also considered. The assumption for the need of random slope was tested using mixtures of chisquares distribution and the results for the test of random effects reduction and model AIC are dispalayed in Table 9 (in appendix A). The result indicated that for all predictors, the more elaborative random effect structures containing random intercepts and random slopes provided a better structure compared to the model with only random intercepts.

Table 5: *Estimated components of the model for the effect of LVEF, heart rhythm and NYHA on DBP resulted from model 3 with random intercepts.*

| LVEF | Smoothed.term | edf | Ref.df | F | p.value |
|---|---|---|---|---|---|
| | s(day)*(LVEF=high) | 6.08 | 6.08 | 4.368 | 0.0002 |
| | s(day)*(LVEF=low) | 4.29 | 4.52 | 10.219 | <0.0001 |
| **Heart rhythm** | | | | | |
| | s(day)*(Heart rhthym=Normal) | 5.84 | 6.58 | 1.99 | 0.055 |
| | s(day)*(Hear trhym=Abnormal) | 6.45 | 7. 23 | 6.75 | <0.0001 |
| **NYHA** | | | | | |
| | s(day)*(NYHA=ClassI) | 1.000 | 1.000 | 29.54 | <0.0001 |
| | s(day)*(NYHA=ClassII) | 2.00 | 2.00 | 0.051 | 0.9500 |
| | s(day)*(NYHA=ClassIII) | 3.53 | 3.53 | 8.64 | <0.0001 |
| | s(day)*(NYHA=ClassIV) | 6.41 | 6.41 | 11.33 | <0.0001 |



Figure 15: *Estimated gam components of model (3) fitted to the DBP. The upper two panels are for LVEF and the lower are for heart rhythm categories.*

Figure 16: *Estimated gam components for model (3) fitted to DBP by level of NYHA.*

## 7.3 Systolic Blood Pressure

A model fitting procedures similar to the one used for diastolic blood pressure was applied for the remaining three biomarkers. For systolic blood pressure, the model comparison results are displayed in Table 10 (in appendix A). In the table the column *-2loglik* is the loglikelihood difference between two models with the same random effects while $-2(ln\lambda_N)$ is the restricted maximum likelihood resulted from mixtures of chi-squares and used to compare models with different random effects. From comparison of model (1) and model (2) with model (3), model (3) was found to be the best fitting model for each variables based on AIC and *-2loglik*. This indicated that along term smoothed function of covariates are needed to describe the systolic trend. We have also extended model (3) by adding linear and quadratic random slopes. These models are represented by model (4) and model (5) in the table respectively. For sex, LVEF and NYHA, model 5 was choosen. For heart rhythm, a model with linear random trend (model (4)

34

was shown as the best model under consideration. Hence, further inference was based on model (5) for sex, LVEF and NYHA while model (4) was used for heart rhythm.

The estimated model components of univariate analysis for the systolic blood pressure from the best fitting models is summarized in Table 6 (in appendix A). A significant non linear smoothed trend is observed for categories of sex and heart rhythm. Having low level of left ventricular ejection fraction (LVEF) was not significant (P=0.0719), indicating that spline is not needed to describe the systolic profiles of patients in this group. Regarding NYHA, a significant smoothed trajectory was obseved for the groups of patients in class II and class IV with (p=0.009 and 0.0005) respectively. The plots of gam components of the model for effect of sex, LVEF and heart rhythm categories on the evolution of systolic blood pressure is given in Figure 17, with effective degrees of freedom on y-axis and x-axis representing time (in days). For male patients the trend shows an initial decreasing followed by increasing phase. On contrary, for female groups, the trend is initially increasing, after reaching the first peak, it evolved with little fluctuation and then the second peak is observed around 150 days. The profile is decreasing for patients with high level of left ventricular ejection fractions (LVEF) while it is increasing for patients with low level of LVEF. For patients in the subgroups of normal heart rhythm, it is initailly decreasing and then increasing for the rest period with some fluctuations. For patients with abnoramal heart rhythm, the estimated trend evolved in a constant fashion except some bumps at about day 100. On the other hand, the the estimated effect along with confidence bands for patients in subgroups of NYHA is depicted in Figure 20 (in appendix B). Non linear trends are evident from the plots.

In conclusion, all results indicated that the longitudinal systolic post discharge profile is better described in a smoothed manner rather than a constant fashion. We have also observed a difference in evolution across the categories of covariates.

Figure 17: *Estimated model terms for additive effect of gender, LVEF and heart rhythm on systolic blood pressure.*

## 7.4   Weight

In Table 11 (appendix A) all models fitted to weight are presented with their respective AIC and loglikelihood. In comparing the common trend, model (1) with the model which assumes different curve for groups, model (3), both AIC and loglikelihood supported the later for all the considered covariates.

To account for subject to subject variability other than at baseline, the linear and quadratic random slopes were considered. These models are given by model (4) and model (5) in the table. However, since model 5 comes with negligible variance for the quadratic random slopes and it did not converged for the sex and heart rhythm covariates, it was not considered. We compared model (3) with model (4) using mixtures of chisquares and model 4 was choosen as optimal to describe the weight profiles for all predictors. The estimated components from model (4) are presented in Table 7 (appendix A). The approximate p-value indicates that the weight curve for the subgroups of patients, i.e male (<0.0001), females (p=<0.0001), patients with low level of LVEF (p=<0.0001) and with normal heart rhythm (p=0.0001) had a smoothed non linear

36

trend. On the other hand, for patients in the first (I), third (III) and fourth (IV) classes of the functional classifications of heart failure (NYHA), the trend had a significant smooth trajectory while for patients in class II, it did not have significant smoothed non linear trends. The resulting plot for estimated components along with the 95% bayesian credible intervals is presented in Figure 18 and in Figure 21 (in appendix B). From both figures it can be observed that the trajectories had some flexible curvatures, except for some subgroups. The coincidence of the confidence limits and the estimated straight line, at the point where the line passes through zero on the vertical axis, is a result of the identifiability constraints applied to the smooth terms.

## 7.5   Heart rate

Different GAMM models were fitted and for all predictors, model (5) was found to to be the best fitting model (Table 12 in appendix A) and adopted for further inference. On the other hand, from comparison of the fit for model (1) and models adjusted for covariates (model 3) the later resulted in smallest AIC and flexible structures for the profiles. The numerical output for model (5) is presented in Table 8 (appendix A). The model output show that there was a signifcant non-linear smoothed trend for both gender and LVEF levels. The results demonstrate that the patients in each categories of sex and LVEF had different evolution trend for heart rate. The assumptions of non linear smoothed trend was not significant for the subgroup with abnormal heart rythm and for the I and III classes of NYHA. This indicated that the long term heart rate profiles for patients in these groups did not have a non linear trajectory rather can be described using other approach, may be linear trend. The gam componets of the model with 95% credible intervals for sex, LVEF and heart rhythm are displayed in Figure 19 and in Figure 22 (in appendix B) for NYHA. From both plots some flexible curvatures are evident with increasing and decreasing trends over time. The baseline predictors, age and NTProBNP were also included in the model in their parametric form and age was found to have a significant effect but not NTproBNP. All in all, flexible relationships of the biomarkers with time were observed from the considerd GAMM models and differences in evolution among patients categorized according to their baseline characteristics. For model diagnostics, the standard conditional residuals against

time was plotted for each outcomes under the choosen optimal models (not presented) and the residuals showed agreat deal of scatter which shows that the considered models approximated the data well.



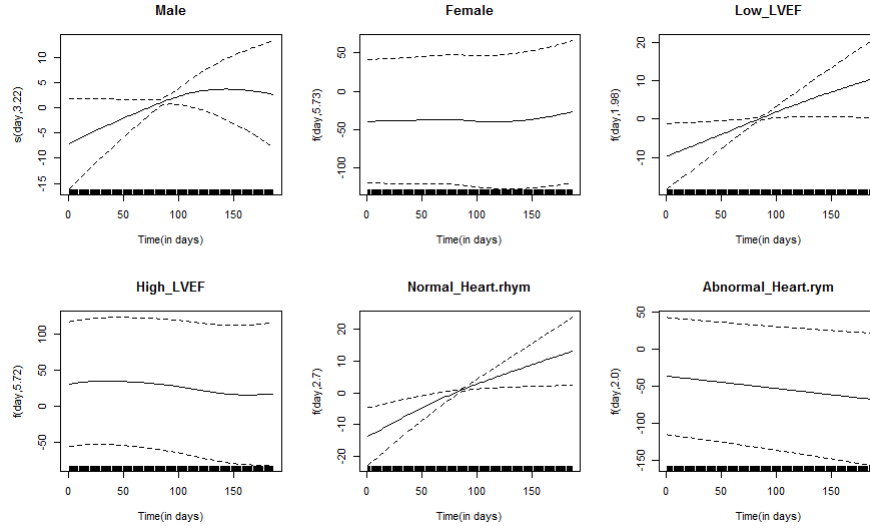Figure 18: *Estimated model terms for additive effect of gender, LVEF and heart rhythm on weight curve obtained from model (4).*
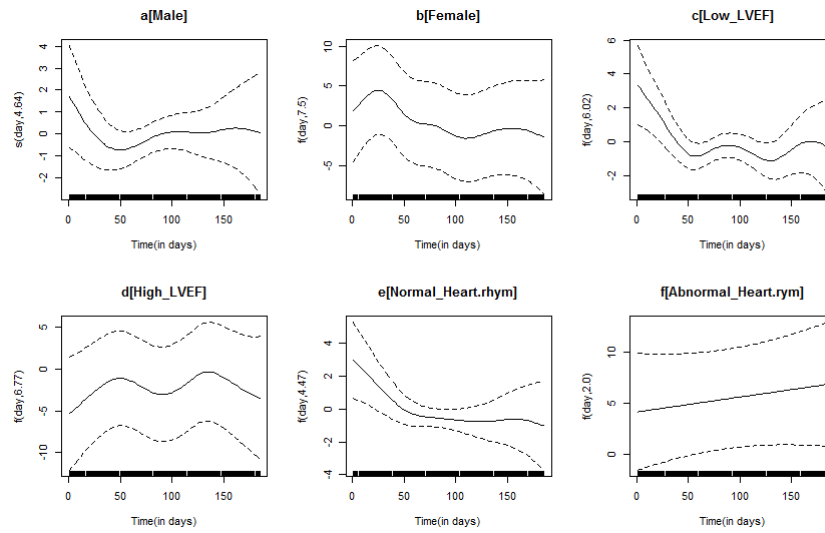


Figure 19: *Estimated model terms for additive effect of gender, LVEF and heart rhythm on heart rate curves obtained from model (5).*

# 8    Discussion and Conclusion

This study based on the analysis of the chronic heart fluire data collected collected in 2008/2010 from different hospitals in Belgium. The objectives of this report are to study the post discharge profiles of four biomarkers measurements collected via telemonitoring instrument after hospitalization and to investigate whether this evolutions are influenced by the patients age, gender, status of left ventricular ejection fraction (LVEF), heart rhythm, measures of cardiac muscle fiber stretch (NTproBNP), and on measures of physical fitness (NYHA). The measurements were taken within a period of six months after hospitalization. All patients did not stay the whole six months. Some patients rehospitalized, some of them dropped from the study. The data has also large number of missing observations due to unknown reason of missingness.

This study has touched upon flexible modelling techniques, with emphasize on penalized spline methodology for modelling the biomarkers profiles. The observed subject specific longitudinal profiles of the four biomarkers showed a curvature like cyclic trend which arise a need to model the time effect in a smoothed fashion in order to better visualize the post discharge profiles. Generalized additive mixed models with p-splines were applied to describe the post discharge profiles in a flexible smoothed way. For all biomarkers, the GAMM resulted in a great deal of flexibility in capturing the relationships of the biomarkers measurement and time. A significant difference in long term evolution was observed for subgroups of patients who where categorized according to their baseline characteristics. The overall and the difference curve showed both increasing and decreasing trajectories for the biomarkers. Patients in certain categories resulted in an unusual trend where the profiles showed a sharp decreasing or increasing trend over time. While others showed both increasing and decreasing trajectories. For instance the diastolic profile for patients whose physical fitness was categorized in the first class according to the NewYork Heart Association (NYHA) functional classification, had a sharp decreasing estimated mean trend for the whole period. Although, prediction is not the objectives of this report, such trends might give a clue for clinicians to take timely intervention. Compared to male patients, the systolic profiles for females showed slight increase within the first month and then

steady until the end. It was aligned with the finding that female patients are more likely to have preserved systolic function (Dumitru and Baker, 2014). In general, we have demonstrated the versatility of the the generalized additive mixed models using penalized spline approach. Clinicians can benefit from this approach to predict patients condition and there by take timely interventions.

## 8.1 Future Reaserch and Recommendations

This report has several limitations. First, the data had missing observations but, we ignored the missingness mechanism based on the assumption of missing at random (MAR) which might not hold true. So future analysis should consider the missingness in the data to varify the validity of the result. Second, we fitted the simple GAMM models where the random effects intered in linear way but in practice semi-parametric models can be fitted by allowing non linear smoothed random effect structures. We carried out univariate analysis by considering the covariate effect at a time. Therefore it is worthwhile to consider the covariates at the same time. Third, we used p-splines with default numbers of knots by assuming equally spaced knots. So it is important to see how sensitive is the fit by varying the number of knots and comparing the result with other spline bases which do not assume equally spaced knots. The selection of the knots (break point) can also be determined in such a way that it is clinically meaningful. We suggest to consider also the clinicians point of view in determining the break points while fitting the model. Last, the four biomarkers were considered separately, so it might be interesting to consider simultaneously using joint models such as multivariate analysis and observe how the evolution changes.

# 9 References

Brumback BA, Ruppert D and Wand MP. (1999). Comment on: Variable selection and function estimation in additive non-parametric regression using a data-based prior, by Thomas SS, Robert K, and Sally W. Journal of the American Statistical Association; **94**: 777-797.

Bui AL and Fonarow GC. (2012). Home monitoring for heart failure management. *J Am Coll Cardiol*; **59**: 97-104.

Chul-Ho Kim *et al.* (2012). A Multivariable Index for Grading Exercise Gas Exchange Severity in Patients with Pulmonary Arterial Hypertension and Heart Failure. *Pulmonary Medicine*; doi:**10**.1155/2012/962-598.

Costa MJ, (2008). Penalized spline models and applications. *University of Warwick,Ph.D. Thesis.* Accessed on 10th July 2014.
Available at: http://wrap.warwick.ac.uk/3654/1/WRAP_THESIS_Costa_2008.pdf.

Craven P and Wahba G. (1979). Smoothing noisy data with spline functions: Estimating the correct degree of smoothing by the method of generalized cross validation.*Numerische Mathematik*;**31**:377-403.

Currie ID and Durban M. (2002). Flexible smoothing with P-splines: A unified approach. *Statistical Modelling* **4**: 333-349.

Dendale P, De Keulenaer G, *et al.* (2011). Effect of a telemonitoring-facilitated collaboration between general practitioner and heart failure clinic on mortality and rehospitalization rates in severe heart failure: the TEMA-HF 1 (TElemonitoring in the MAnagement of Heart Failure) study. *European Journal of Heart Failure*; **14**: 333-340.

Dickstein K, Cohen-Solal A, *et al.* (2008). ESC Committee for Practice Guidelines (CPG). ESC guidelines for the diagnosis and treatment of acute and chronic heart failure 2008: the Task Force for the diagnosis and treatment of acute and chronic heart failure 2008 of the European Society of Cardiology. *Eur J Heart Fail*; **10**: 933-989.

Dumitru I and Baker MM. (2014). Heart Failure. *Medscape.* Accessed on 17th August 2014. Available at: http://emedicine.medscape.com/article/163062-overview#a0156.

Durbán M, Harezlak J, *et al.* (2005). Simple Fitting of Subject Specific Curves for Longitudinal Data. *Statistics in Medicine*; **24**: 1153-1162.

Eilers PHC and Marx BD. (1996). Flexible smoothing with B-splines and penalties. *Statistical Science*; **11**: 89-121.

Fitzmaurice GM, Laird NM, and Ware JH. (2009). *Longitudinal Data Analysis: Handbooks of Modern Statistical Methods.* Boca Raton, Florida: Chapman & Hall/CRC.

Friedman JH and Silverman BW. (1989). *Flexible parsimonious smoothing and additive modeling. Technometrics*; **31**: 321.

Hastie TJ and Tibshirani RJ. (1990). *Generalized Additive Models.* Boca Raton, Florida: Chapman & Hall.

Ho KK, Pinsky JL, Kannel WB and Levy D. (1993). The epidemiology of heart failure: the Framingham Study. *J Am Coll Cardiol*; **22** Suppl 4: A6-13.

Jencks SF, Williams MV and Coleman EA. (2009). Rehospitalizations among patients in the Medicare fee-for-service program. *N Engl J Med*; **360**: 1418-1428.

Lee DJ. (2010). Smoothing mixed models for spatial and spatio-temporal data.*Universidad Carlos III de Madrid, Ph.D. Thesis.* Accessed on 12th August 2014. Available at: http://e-archivo.uc3m.es/bitstream/handle/10016/9364/Tesis_DaeJinLee.pdf?sequence=1

Lee DJ and Durbán M. (2009). Smooth-CAR mixed models for spatial count data. *Computational Statistics and Data Analysis*; **53**: 2968-2977.

Lin X and Zhang D. (1999). Inference in generalized additive mixed models using smoothing splines. *Journal of the Royal Statistical Society*; Series B **61**: 381-400.

Maringwa JT, Geys H, *et al.* (2008c). Application of semi-parametric mixed models and simultaneous confidence bounds in a cardiovascular safety experiment with longitudinal data. *Journal of Biopharmaceutical Statistics*, **18**(6), 1043-1062.

Marx BD and Eilers PHC. (1998). Direct generalized additive modeling with penalized likelihood. *Computational Statistics and Data Analysis*; **28**: 193-209.

Njagi NE, Molenberghs G, Rizopoulus D, *et al.* (2013a). A flexible joint-modeling framework for longitudinal and time-to-event data with overdispersion. *Statistical Methods in Medical Research*. Published online before print July 18, 2013, doi: 10.1177/0962280213495994.

Njagi NE, Rizopoulos D, Molenberghs G, *et al.* (2013b). A Joint Survival-Longitudinal Modelling Approach for the Dynamic Prediction of Rehospitalization in Telemonitored Chronic Heart Failure Patients. *Statistical Modelling* ; **13(3)**: 179-198.

Njagi NE. (2009). Longitudinal analysis of fast fluorescent induction in plants to study the effects of non-photosynthetic oxidative stresses on the photosynthetic process. Unpublished Masters thesis. Hasselt University, BELGIUM.

O'Sullivan F. (1986). A statistical perspective on ill-posed inverse problems (with discussion).*Statistical Science*; **1**: 505-527.

Pinheiro JC and Bates DM. (2000). *Mixed effects models in S and S-Plus.* New York: Springer.

Ruppert D, Wand MP and Carroll RJ. (2003) *Semiparametric Regression.* Cambridge: Cambridge University Press.

Verbeke G and Molenberghs G. (2000). *Linear Mixed Models for Longitudinal Data.* New York: Springer.

Wood SN. (2008). Fast stable direct fitting and smoothness selection for generalized additive models. *Journal of the Royal Statistical Society*; Series B **70**: 495-518.

Wood SN. (2006). *Generalized Additive Models An Introduction with R.* Boca Raton, Florida: Chapman & Hall/CRC.

# 10 APPENDIX

## 10.1 Appendix A:Tables

Table 6: *Estimated gam components of the model for additive effects of Sex, LVEF, heart rhythm and NYHA on SBP resulted from model (5) for sex and LVEF and model (4) for heart rhythm.*

| Predictor | Smoothed.term | edf | Ref.df | F | p.value |
|---|---|---|---|---|---|
| **Sex** | | | | | |
| | s(day)*(Sex=M) | 5.011 | 5.031 | 5.873 | <0.0001 |
| | s(day)*(Sex=F) | 6.685 | 6.687 | 6.099 | <0.0001 |
| **LVEF** | | | | | |
| | s(day)*(LVEF=high) | 2.000 | 2.001 | 6.584 | 0.0014 |
| | s(day)*(LVEF=low) | 3.342 | 3.352 | 2.263 | 0.0719 |
| **Heartrhythm** | | | | | |
| | s(day)*(Heart rhythm=Normal) | 3.046 | 3.046 | 5.377 | 0.0010 |
| | s(day)*(Heart rhythm=Abnormal) | 6.462 | 6.470 | 3.457 | 0.0016 |
| **NYHA** | | | | | |
| | s(day)*(NYHA=ClassI) | 1.001 | 1.001 | 0.860 | 0.3536 |
| | s(day)*(NYHA=ClassII) | 4.742 | 4.742 | 3.133 | 0.0094 |
| | s(day)*(NYHA=ClassIII) | 4.171 | 4.171 | 1.067 | 0.1500 |
| | s(day)*(NYHA=ClassIV) | 5.882 | 5.882 | 4.054 | 0.0005 |

Table 7: *Estimated gam components of the model for additive effects of Sex, LVEF, heart rhythm and NYHA on Weight from model (4).*

| Predictor | Smoothed.term | edf | Ref.df | F | p.value |
|:---:|:---|---:|---:|---:|---:|
| **Sex** | | | | | |
| | s(day)*(Sex=M) | 3.216 | 3.217 | 8.346 | <0.0001 |
| | s(day)*(Sex=F) | 5.733 | 5.735 | 5.454 | <0.0001 |
| **LVEF** | | | | | |
| | s(day)*(LVEF=high) | 1.983 | 1.982 | 2.860 | 0.0580 |
| | s(day)*(LVEF=low) | 5.716 | 5.720 | 7.806 | <0.0001 |
| **Heart rhythm** | | | | | |
| | s(day)*(Heart rhythm=Normal) | 2.696 | 2.696 | 7.27 | 0.0001 |
| | s(day)*(Heart rhythm=Abnormal) | 2.003 | 2.003 | 2.31 | 0.0992 |
| **NYHA** | | | | | |
| | s(day)*(NYHA=ClassI) | 0.999 | 0.999 | 0.613 | <0.0001 |
| | s(day)*(NYHA=ClassII) | 2.109 | 2.109 | 0.180 | 0.846 |
| | s(day)*(NYHA=ClassIII) | 6.793 | 6.793 | 3.755 | 0.0005 |
| | s(day)*(NYHA=ClassIV) | 7.200 | 7.200 | 2.375 | 0.0190 |

Table 8: *Estimated gam components of the model for additive effects of Sex, LVEF, heart rhythm and NYHA on heart rate from model (5).*

| Predictor | Smoothed.term | edf | Ref.df | F | p.value |
|:---:|:---|---:|---:|---:|---:|
| **Sex** | | | | | |
| | s(day)*(Sex=M) | 4.637 | 4.639 | 5.360 | 0.0001 |
| | s(day)*(Sex=F) | 7.496 | 7.496 | 4.835 | <0.0001 |
| **LVEF** | | | | | |
| | s(day)*(LVEF=high) | 6.772 | 6.772 | 6.148 | <0.0001 |
| | s(day)*(LVEF=low) | 6.017 | 6.017 | 7.343 | <0.0001 |
| **Heartrhythm** | | | | | |
| | s(day)*(Heart rhthym=Normal) | 4.474 | 4.474 | 4.219 | 0.0014 |
| | s(day)*(Heart rhythm=Abnormal) | 2.002 | 2.002 | 2.805 | 0.0605 |
| **NYHA** | | | | | |
| | s(day)*(NYHA=ClassI) | 4.351 | 4.351 | 2.282 | 0.0529 |
| | s(day)*(NYHA=ClassII) | 7.478 | 7.478 | 19.493 | <0.0001 |
| | s(day)*(NYHA=ClassIII) | 2.001 | 2.001 | 2.192 | 0.1117 |
| | s(day)*(NYHA=ClassIV) | 2.010 | 2.010 | 3.145 | 0.0429 |

# Models Comparison

## I.Diastolic blood pressure

Table 9: *Test for reduction of random effects and model AIC for DBP .*

|  | Model | AIC | logLik | -2loglik | p-value |  |
|---|---|---|---|---|---|---|
| **LVEF** | Model (1) | 66971.99 | -33478.00 |  |  |  |
|  | Model (2) | 66973.02 | -33480.38 | 0.12 | 0.7300 |  |
|  |  | AIC | logLik |  | -2ln($\lambda_N$) | p-value |
|  | Model(3)_RI | 65825.71 | -32901.85 |  |  |  |
|  | Model (3)_RS | 65581.10 | -32776.55 |  | 250.60 | <0.0001 |
| **Heart rhythm** | Model (1) | 66971.99 | -33480.57 |  |  |  |
|  | Model (2) | 66972.94 | -33480.47 | 0.21 | 0.6500 |  |
|  |  | AIC | logLik |  | -2ln($\lambda_N$) | p-value |
|  | Model(3)_RI | 65826.35 | -32902.17 |  |  |  |
|  | Model(3)_RS | 65604.17 | -32788.08 |  | 228.1808 | <0.0001 |
| **NYHA** | Model (1) | 66971.99 | -33478.00 |  |  |  |
|  | Model (2) | 66974.50 | -33479.25 | 2.65 | 0.4494 |  |
|  |  | AIC | logLik |  | -2ln($\lambda_N$) | p-value |
|  | Model(3)_RI | 66895.75 | -33433.88 |  |  |  |
|  | Model(3)_RS | 65826.68 | -32896.34 | 1075.07 | < 0.0001 |  |

## II.Systolic blood pressure

Table 10: *Test for reduction of random effects and model AIC for SBP.*

| Sex | Model | AIC | logLik | -2loglik | -2ln($\lambda_N$) | p-value |
|---|---|---|---|---|---|---|
| | Model 1 | 74657.57 | -37323.78 | | | |
| | Model 2 | 74651.68 | -37319.84 | 7.88 | | 0.005 |
| | Model 3 | 74621.53 | -37319.84 | 34.15 | | < .0001 |
| | Model 4 | 72987.86 | -36482.93 | | 1639.67 | < .0001 |
| | Model 5 | 72637.91 | -36304.95 | | 355.95 | < .0001 |
| **LVEF** | Model 1 | 74657.57 | -37323.78 | | | |
| | Model 2 | 74649.57 | -37318.78 | 10.00 | | 0.0016 |
| | Model 3 | 74627.26 | -37305.63 | 26.31 | | < .0001 |
| | Model 4 | 72987.88 | -36482.94 | | 1645.38 | < .0001 |
| | Model 5 | 72649.44 | -36310.72 | | 344.44 | < .0001 |
| **Heart rhythm** | Model 1 | 74657.57 | -37323.78 | | | |
| | Model 2 | 74659.43 | -37323.71 | 0.14 | | 0.7066 |
| | Model 3 | 74643.45 | -37313.73 | 19.97 | | < .0001 |
| | Model 4 | 72986.92 | -36482.46 | | 1662.53 | < .0001 |
| | Model 5 | | | | | not converged |
| **NYHA** | Model 1 | 74657.57 | -37323.78 | | | |
| | Model 2 | 74661.15 | -37322.57 | 2.42 | | 0.4899 |
| | Model 3 | 74297.66 | -37134.83 | 375.49 | | < .0001 |
| | Model 4 | 72975.68 | -36470.84 | | 1327.98 | < .0001 |
| | Model 5 | 72636.61 | -36298.30 | | 345.07 | < .0001 |

## III.Weight

Table 11: *Test for reduction of random effects and model AIC for weight.*

| Sex | Model | AIC | logLik | -2loglik | -2ln($\lambda_N$) | p-value |
|---|---|---|---|---|---|---|
| | Model 1 | 91782.73 | -45886.36 | | | |
| | Model 2 | 91783.75 | -45885.87 | 0.98 | | 0.3221 |
| | Model 3 | 91766.72 | -45875.36 | 21.02 | | < .0001 |
| | Model 4 | 89839.78 | -44908.89 | | 1932.95 | < .0001 |
| **LVEF** | Model 1 | 91782.73 | -45886.36 | | | |
| | Model 2 | 91784.19 | -45886.10 | 0.53 | 0.47 | |
| | Model 3 | 91761.69 | -45872.85 | 26.50 | | < .0001 |
| | Model 4 | 89826.90 | -44902.45 | | 1940.79 | < .0001 |
| **Heart rhythm** | Model 1 | 91782.73 | -45886.36 | | | |
| | Model 2 | 91783.50 | -45885.75 | 1.23 | | 0.27 |
| | Model 3 | 91727.84 | -45855.92 | 59.66 | | < .0001 |
| | Model 4 | 89858.17 | -44918.08 | | 1875.67 | < .0001 |
| **NYHA** | Model 1 | 91782.73 | -45886.36 | | | |
| | Model 2 | 91788.02 | -45886.01 | 0.7034 | | 0.8724 |
| | Model 3 | 91618.17 | -45795.09 | 181.85 | | < .0001 |
| | Model 4 | 89862.75 | -44914.37 | | 1761.42 | < .0001 |

## IV.Heart rate

Table 12: *Test for reduction of random effects and model AIC for heart rate.*

| Sex | Model | AIC | logLik | -2loglik | -2ln($\lambda_N$) | p-value |
|---|---|---|---|---|---|---|
| | Model 1 | 67590.93 | -33790.47 | | | |
| | Model 2 | 67592.63 | -33790.31 | 0.31 | | 0.5811 |
| | Model 3 | 67464.50 | -33724.25 | 132.13 | | < .0001 |
| | Model 4 | 65738.35 | -32858.17 | | 1732.15 | < .0001 |
| | Model 5 | 65267.93 | -32619.96 | | 476.42 | < .0001 |
| **LVEF** | Model 1 | 67590.93 | -33790.47 | | | |
| | Model 2 | 67592.28 | -33790.14 | 0.65 | | 0.4207 |
| | Model 3 | 67584.59 | -33784.30 | 11.69 | | 0.0029 |
| | Model 4 | 65756.20 | -32867.10 | | 1834.39 | < .0001 |
| | Model 5 | 65267.51 | -32619.75 | | 494.69 | < .0001 |
| **Heart rhythm** | Model 1 | 67590.93 | -33790.47 | | | |
| | Model 2 | 67587.93 | -33787.975 | 4.99 | | 0.0254 |
| | Model 3 | 67555.61 | -33769.80 | 36.32 | | < .0001 |
| | Model 4 | 65778.73 | -32878.36 | | 1782.88 | < .0001 |
| | Model 5 | 65288.79 | -32630.39 | | 495.93 | < .0001 |
| **NYHA** | Model 1 | 67590.93 | -33790.47 | | | |
| | Model 2 | 67584.96 | -33783.48 | 13.97 | | 0.2400 |
| | Model 3 | 67414.12 | -33693.06 | 181.08 | | < .0001 |
| | Model 4 | 65676.14 | -32821.07 | | 1743.98 | < .0001 |
| | Model 5 | 65166.92 | -32563.46 | | 515.2207 | < .0001 |

## 10.2    Appendix B:Figures



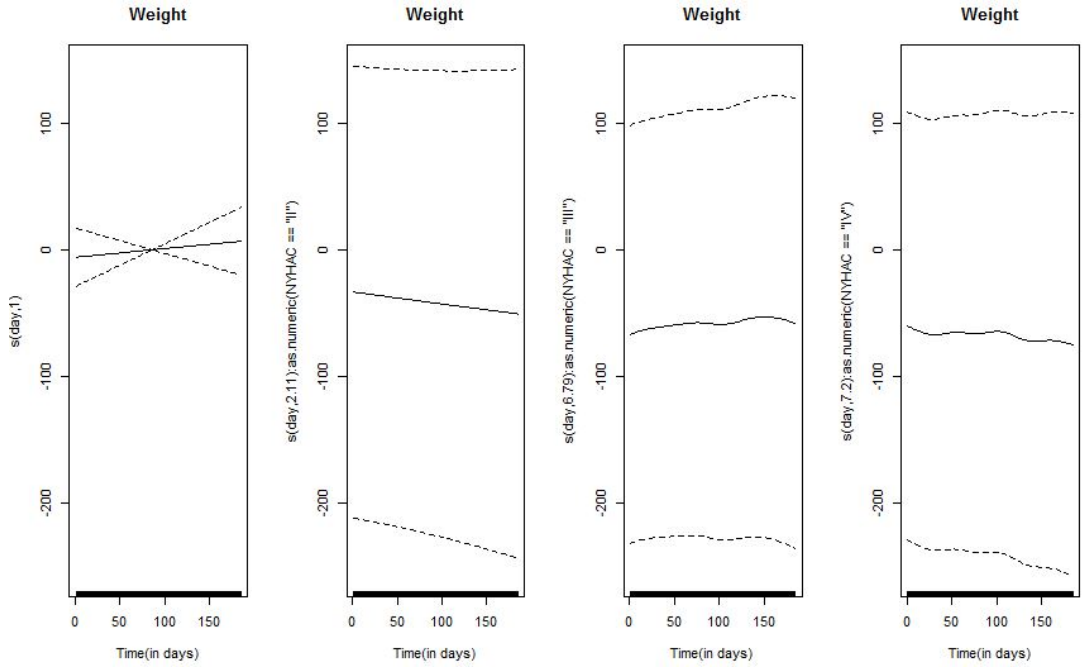Figure 20: *Estimated SBP curve by NYHA from model (5) with 95% bayesian credible intervals.*



Figure 21: *Estimated Weight curves by NYHA from model (4) with 95% bayesian credible intervals.*
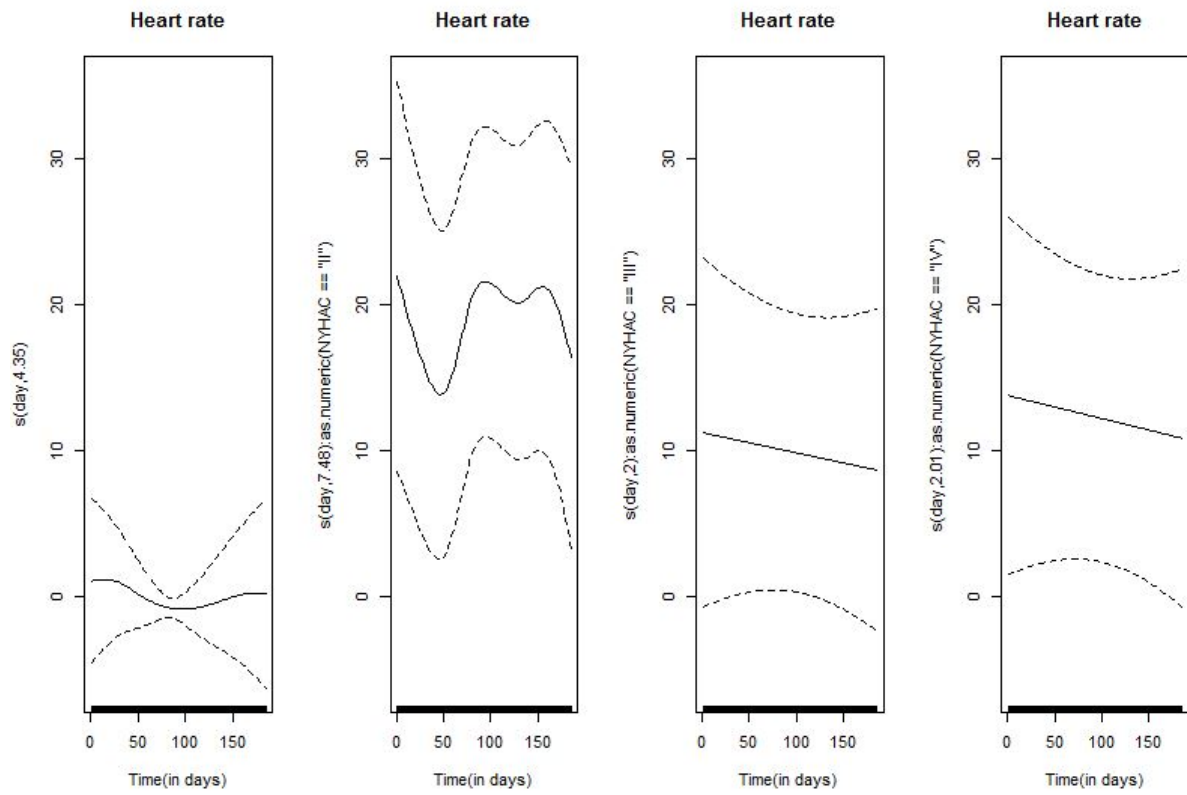
Figure 22: *Estimated heart rate curves by NYHA with 95% bayesian credible intervals from model (5).*

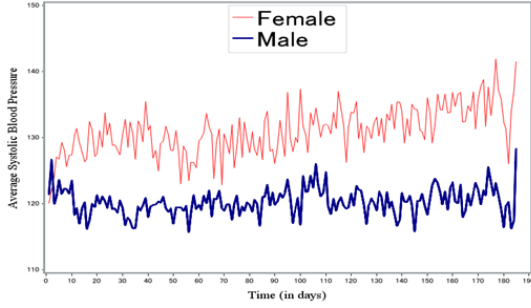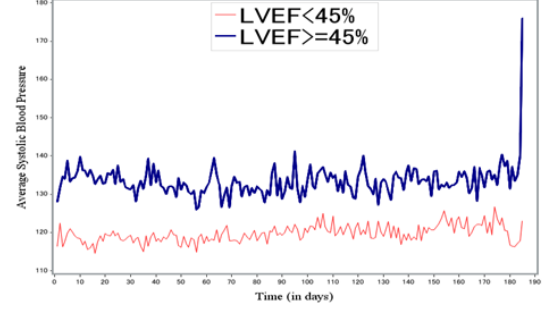# I. Average evolution of Systolic blood pressure



(e)



Figure 23: *Average evolution of systolic blood pressure by levels of covariates.*

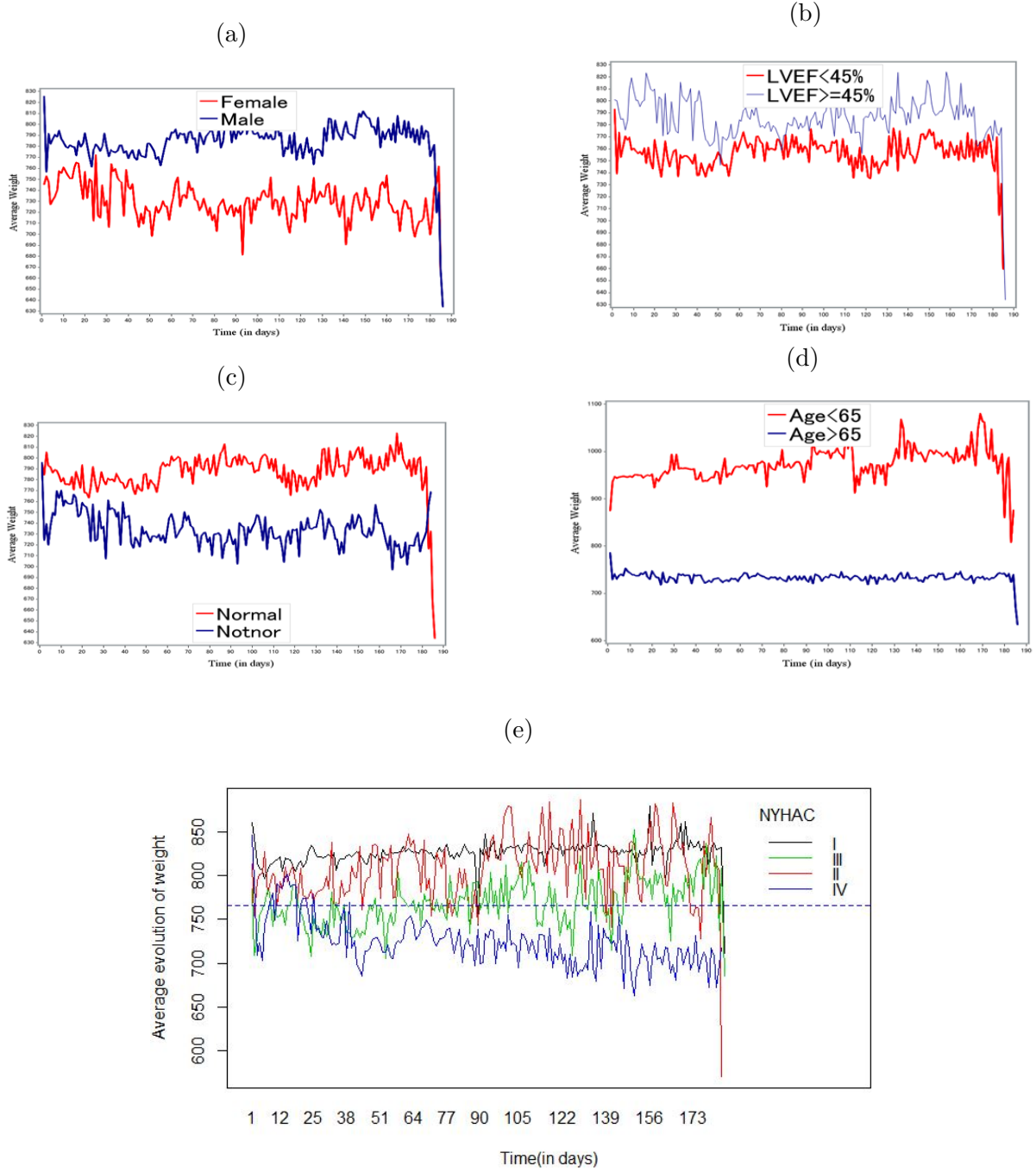**II. Average evolution of Weight**

(a)

(b)

(c)

(d)

(e)



Figure 24: *Average evolution of Weight by gender(a) ,LVEF (b), heart rhythm (c),age (d) and NYHA (e)*

# III. Average evolution of heart rate
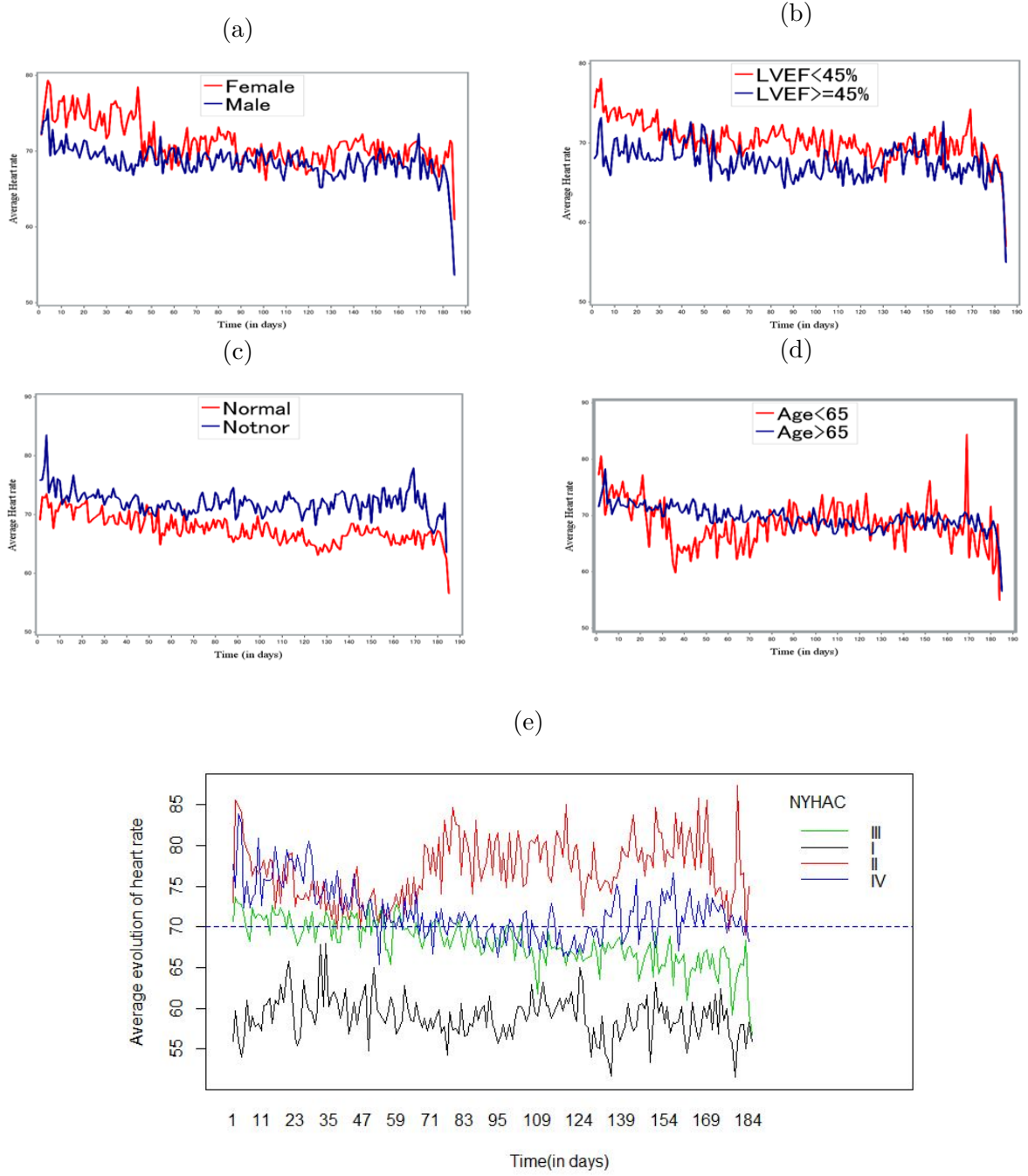
(a)

(b)

(c)

(d)

(e)



Figure 25: *Average evolution of heart rate by gender(a) ,LVEF (b), heart rhythm (c),age (d) and NYHA (e)*
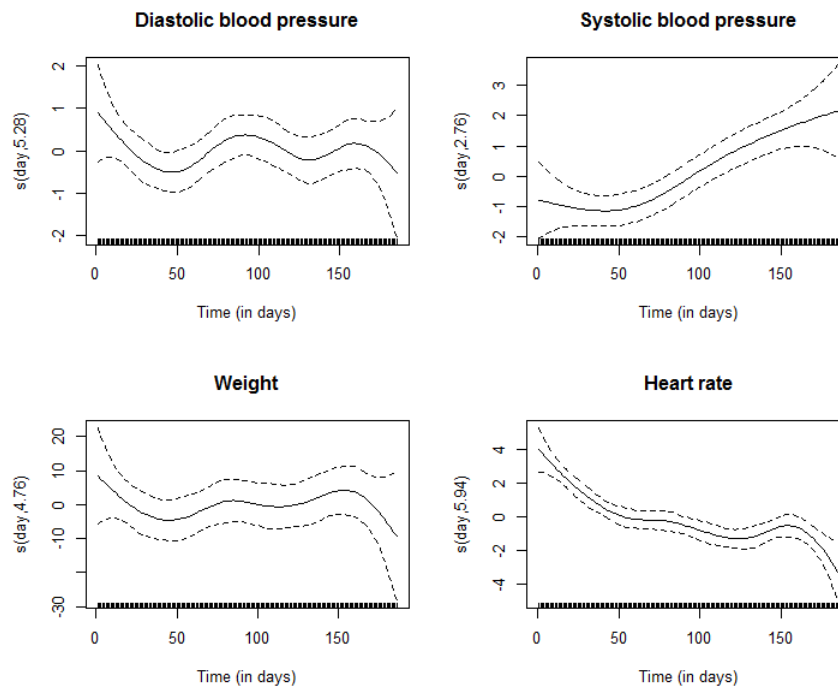
Figure 26: *Over all estimated curves of DBP, SBP, Weight and Heart rate from GAM (without random effects and with only time in the model)*

## 10.3    Appendix E:Codes

### 10.3.1    SAS-code

```
****************************PLOTS*****************
*individual profile plots*;
*****************************************dbp*******************;
proc sort data=thesi1;by ptid time ;run;
goptions reset=all i=join;
axis1 w=3  label=(h=2  minor=none
 font='times new roman' "Daysafterhospitalization")
order=(1 to 186 by 10);
axis2 w=3  label=(h=2 A=90
font='times new roman' 'Diastolic blood pressure  ')
order=(40 to 134 by 5);
```

```
footnote font='times new roman' h=2
'Figure 2: Individual profile plots for diastolic blood pressure';
proc gplot data=thesi1;plot dbp*day=PTID/ haxis=axis1 vaxis=axis2
nolegend;run;quit;


** sample codes for Average evolution by levels of covariates**
/* average DBP plot by SEX*/
proc sort data=thesis;by  sexc day ;
proc means data=thesis;
var dbp  sbp  weight heartrate;by sexc day;
output out=meanNYHA ;run;
data meanNYHA1; set meanNYHA; if _STAT_='MEAN' then output; run;
goptions reset=all ;
symbol1  c=red w=1  i=j;
symbol2  c=darkblue W=1  i=j ;
axis1 W=2 color=DEEPblack label=(h=2
font='times new roman' 'Time (in days)')minor=none
order=( 1 to 200 by 10);;
axis2 W=2 color=DEEPblack label=(h=2 A=90
font='times new roman' 'Average Diastolic Blood Pressure ')
minor=none order=( 60 to 90 by 5);
footnote font='times new roman' h=2
'Figure SBPa: Average profile plots by gender' ;
LEGEND LABEL=NONE POSITION=(TOP  INSIDE) VALUE=(H=4)
DOWN =2 FRAME;
proc gplot data=meanNYHA1;
plot DBP*day=sexc/haxis=axis1 vaxis=axis2 LEGEND=LEGEND;
run;
```

56

```
quit;
```

## 10.3.2  R-code

```
#Codes used to fit GAMM#
rm(list = ls());ls()
#Statistical analysis#;
#GAMMs,HEARTFLUIRE DATA#;#read data#
memory.limit(size = NA)
memory.limit(size = 4000)
#HeartFluire=na.omit(HeartFluire)#
HeartFluire=read.table("C:\\Users\\abdi\\Documents\\SECONDYR
\\sem-2nd\\thesis\\Thesis.csv", sep=",",header=T, na.string="*")
HeartFluire$heartrym<-factor(HeartFluire$heartrym)
HeartFluire$sex<-factor(HeartFluire$sex)
HeartFluire$LVEF <-as.numeric(HeartFluire$LVEF)
HeartFluire$NYHACLASS=as.numeric(HeartFluire$NYHA)
HeartFluire$NYHACLASS=factor(HeartFluire$NYHACLASS)
str(HeartFluire)
#Diastolic MODELS#
#SEX
#over alltrend
model1<-gamm(dbp~s(day,bs="ps"),
      data=HeartFluire,random=list(ptID=~1),control=lmc)
#constant effect of sex
model2<-gamm(dbp~s(day,bs="ps")+sexc,
       data=HeartFluire,random=list(ptID=~1),control=lmc)
#difference curve with random intercept
```

```r
model3<-gamm(dbp~s(day,bs="ps")+s(day,bs="ps",by=as.numeric(sexc=="Female")),
      data=HeartFluire,random=list(ptID=~1),control=lmc)
#difference curve with random intercept and slopes
model4<-gamm(dbp~s(day,bs="ps")+s(day,bs="ps",by=as.numeric(sexc=="Female")),
data=HeartFluire,random=list(ptID=~1,ptID=~day),control=lmc)
#difference curve with random intercept, linear and quadratic random slopes
model5<-gamm(dbp~s(day,bs="ps")+s(day,bs="ps",by=as.numeric(sexc=="Female")),
data=HeartFluire,random=list(ptID=~1,ptID=~day+I(day^2)),control=lmc)


#LVEF
#over alltrend
model1<-gamm(dbp~s(day,bs="ps"),
    data=HeartFluire,random=list(ptID=~1),control=lmc)
#constant effect of LVEF
model2<-gamm(dbp~s(day,bs="ps")+LVEFC,
     data=HeartFluire,random=list(ptID=~1),control=lmc)
#difference curve with random intercept
model3<-gamm(dbp~s(day,bs="ps")+s(day,bs="ps",by=as.numeric(LVEFC=="HIGH")),
data=HeartFluire,random=list(ptID=~1),control=lmc)
 #difference curve with random intercept and slopes
 model4<-gamm(dbp~s(day,bs="ps")+s(day,bs="ps",by=as.numeric(LVEFC=="HIGH")),
data=HeartFluire,random=list(ptID=~1,ptID=~day),control=lmc)
 #difference curve with random intercept and slopes+random quadratic slopes
 model5<-gamm(dbp~s(day,bs="ps")+s(day,bs="ps",by=as.numeric(LVEFC=="HIGH")),
 data=HeartFluire,random=list(ptID=~1,ptID=~day+I(day^2)),control=lmc)
#Heart rhythm
model1<-gamm(dbp~s(day,bs="ps"),
        data=HeartFluire,random=list(ptID=~1),control=lmc)
```

```
model2<-gamm(dbp~s(day,bs="ps")+heartrymc,
         data=HeartFluire,random=list(ptID=~1),control=lmc)
model3<-gamm(dbp~s(day,bs="ps")+s(day,bs="ps",by=as.numeric(heartrymc=="Not normal")),
          data=HeartFluire,random=list(ptID=~1),control=lmc)
model4<-gamm(dbp~s(day,bs="ps")+s(day,bs="ps",by=as.numeric(heartrymc=="Not normal")),
         data=HeartFluire,random=list(ptID=~1,ptID=~day),control=lmc)
model5<-gamm(dbp~s(day,bs="ps")+s(day,bs="ps",by=as.numeric(heartrymc=="Not normal")),
       data=HeartFluire,random=list(ptID=~1,ptID=~day+I(day^2)),control=lmc)
#NYHA
model1<-gamm(dbp~s(day,bs="ps"),
       data=HeartFluire,random=list(ptID=~1),control=lmc)
model2<-gamm(dbp~s(day,bs="ps")+NYHACLASS,
       data=HeartFluire,random=list(ptID=~1),control=lmc)
model3<-gamm(dbp~s(day,bs="ps")+
           s(day,bs="ps",by=as.numeric(NYHACLASS=="II"))+
           s(day,bs="ps",by=as.numeric(NYHACLASS=="III"))+
           s(day,bs="ps",by=as.numeric(NYHACLASS=="IV")),
           data=HeartFluire,random=list(ptID=~1),control=lmc)
model4<-gamm(dbp~s(day,bs="ps")+
         s(day,bs="ps",by=as.numeric(NYHACLASS=="II"))+
         s(day,bs="ps",by=as.numeric(NYHACLASS=="III"))+
         s(day,bs="ps",by=as.numeric(NYHACLASS=="IV")),
          data=HeartFluire,random=list(ptID=~1,ptID=~day),control=lmc)
model5<-gamm(dbp~s(day,bs="ps")+
       s(day,bs="ps",by=as.numeric(NYHACLASS=="II"))+
       s(day,bs="ps",by=as.numeric(NYHACLASS=="III"))+
       s(day,bs="ps",by=as.numeric(NYHACLASS=="IV")),
       data=HeartFluire,random=list(ptID=~1,ptID=~day+I(day^2)),control=lmc)
```

```
#Example codes to extract summary from the fitted model

#GAM COMPONENT

summary(model5$gam)

#LME COMPONENT

summary(model5$lme)


#Plotting GAM componnets

    par(mfrow=c(2,2))

 plot(model5$gam,xlab="Time(in days)",main="a[Male]",select=1,scale=0)

  plot(model5$gam,xlab="Time(indays)",main="b[Female]",select=2,scale=0,

    ylab="f(day,5.73)")##for ploting


#Plotfor residual diagnostics

    x<-model5$lme

    plot(x,resid(.,type="p")~day)
```

# Auteursrechtelijke overeenkomst

Ik/wij verlenen het wereldwijde auteursrecht voor de ingediende eindverhandeling:
**Flexible Modelling of Postdischarge Profiles in Telemonitored Chronic Heart Failure Patients**

Richting: **Master of Statistics-Biostatistics**
Jaar: **2014**

in alle mogelijke mediaformaten, - bestaande en in de toekomst te ontwikkelen - , aan de Universiteit Hasselt.

Niet tegenstaand deze toekenning van het auteursrecht aan de Universiteit Hasselt
behoud ik als auteur het recht om de eindverhandeling, - in zijn geheel of gedeeltelijk -,
vrij te reproduceren, (her)publiceren of distribueren zonder de toelating te moeten
verkrijgen van de Universiteit Hasselt.

Ik bevestig dat de eindverhandeling mijn origineel werk is, en dat ik het recht heb om de rechten te verlenen die in deze overeenkomst worden beschreven. Ik verklaar tevens dat de eindverhandeling, naar mijn weten, het auteursrecht van anderen niet overtreedt.

Ik verklaar tevens dat ik voor het materiaal in de eindverhandeling dat beschermd wordt door het auteursrecht, de nodige toelatingen heb verkregen zodat ik deze ook aan de Universiteit Hasselt kan overdragen en dat dit duidelijk in de tekst en inhoud van de eindverhandeling werd genotificeerd.

Universiteit Hasselt zal mij als auteur(s) van de eindverhandeling identificeren en zal geen wijzigingen aanbrengen aan de eindverhandeling, uitgezonderd deze toegelaten door deze overeenkomst.

Voor akkoord,

**Wele, Abduljewad Nuru**

Datum: **10/09/2014**