# Flexible semiparametric analysis of longitudinal genetic studies by reduced rank smoothing

**Yuanjia Wang**[*],
Columbia University, New York, USA

**Chiahui Huang**,
Columbia University, New York, USA

**Yixin Fang**,
Georgia State University, Atlanta, USA

**Qiong Yang**, and
Boston University, Boston, USA

**Runze Li**
The Pennsylvania State University at University Park, University Park, USA

## Abstract

In family-based longitudinal genetic studies, investigators collect repeated measurements on a trait that changes with time along with genetic markers. Since repeated measurements are nested within subjects and subjects are nested within families, both the subject-level and measurement-level correlations must be taken into account in the statistical analysis to achieve more accurate estimation. In such studies, the primary interests include to test for quantitative trait locus (QTL) effect, and to estimate age-specific QTL effect and residual polygenic heritability function. We propose flexible semiparametric models along with their statistical estimation and hypothesis testing procedures for longitudinal genetic designs. We employ penalized splines to estimate nonparametric functions in the models. We find that misspecifying the baseline function or the genetic effect function in a parametric analysis may lead to substantially inflated or highly conservative type I error rate on testing and large mean squared error on estimation. We apply the proposed approaches to examine age-specific effects of genetic variants reported in a recent genome-wide association study of blood pressure collected in the Framingham Heart Study.

### Keywords

Genome-wide association study; Penalized splines; Quantitative trait locus

## 1 Introduction

For quantitative traits that change with age, such as blood pressure and cholesterol level, longitudinal genetic studies can offer valuable opportunity to detect genes that have a time-varying effect and examine how genes affect developmental features of these traits. One example of a longitudinal genetic study is the Framingham Heart Study (FHS) (Dawber et al. 1951), a large ongoing prospective study of risk factors for cardiovascular disease (CVD) originated in 1948. Since its initiation, the study has produced many major discoveries that

[*]Address for Correspondence: Yuanjia Wang, Department of Biostatistics, Mailman School of Public Health, Columbia University, New York, NY 10032, USA, yuanjia.wang@columbia.edu.

have helped scientists understand the development and progression of heart disease and its risk factors. In the FHS, repeated measurements are collected on subjects' clinical characteristics such as cholesterol level, blood pressure, and blood glucose. To understand genetic underpinning of risk factors for CVD, dense single-nucleotide polymorphism (SNP) genotyping was performed using approximately 550,000 SNPs in nearly ten thousand individuals from three generation families in the FHS. The FHS provides an opportunity to discover not only genes affecting mean value of a risk factor, but also the ones that affect time-varying features such as rate of change over time in a trait.

Theories and evidences for genetic factors controlling time-varying developmental features of a phenotype are noted in plant, animal and human genetics literature. For example, complex biological organisms such as plants and animals have evolved through accumulation of mutations in genes that control the developmental processes leading to their mature forms (Rice 2002, Raff 2000, He et al. 2010). Rice (2002) described general population genetic models which construct differential equations to relate developmental features of traits to QTLs. From an evolutionary and developmental biology perspective, Raff et al. (2000) discussed mechanisms of regulatory genes controlling developmental features of complex organisms. Zhao et al. (2004) mapped genes controlling rice plant growth which suggested plants with certain genes would manifest faster growth. In human genetics, Province and Rao (1985) observed temporal trends for familial aggregation and heritability of systolic blood pressure in a Japanese-American family study, and Jarvik et al. (1997) demonstrated age-dependent effect of the apo-E genotype on lipid levels.

Despite these evidence, however, interactions among gene and age or age-dependent genetic effect is routinely ignored in genetic analysis (Lasky-Su et al. 2008). One disadvantage of such practice is that it may make discovery of individual genes with moderate effects more difficult due to loss of power (Shi and Rao 2008). Another limitation is that it may contribute to inconsistent replication of genetic association findings (Lasky-Su et al. 2008). For example, when there exists gene-age interaction, subjects in a replication sample may be in a different age range than the initial study sample so that the replication study may fail to discover a gene that has an effect in the original study age range. Van Steen et al. (2005) proposed screening and testing algorithms for replication within a single set of family data.

A naive way of analyzing longitudinal genetic data is to perform a set of genetic analyses at each time point separately (Atwood et al. 2002). However, this approach ignores rich information in the longitudinal structure and may not detect genes affecting time-varying features of a trait. Strauch et al. (2003) reviewed several two-step methods: the first step is either to take the average of trait measurements on a subject or to fit a longitudinal model without consideration of genetic markers or family structures; the second step is to perform genetic analysis on one or more summary statistics derived from the first step. This method may be improved by a joint approach that fits longitudinal and genetic parameters simultaneously. Shi and Rao (2008) and Zhang and Zhong (2006) used a parametric function such as exponential or Gaussian to accommodate time-varying genetic effect in linkage studies. Wu, Ma and Casella (2007) summarized parametric methods on functional mapping of genes with time-varying effect through a Gaussian mixture model with controlled population. Shi and Rao (2008) showed that ignoring temporal trends in genetic effects can reduce power substantially. While the major advantage of parametric models is its parsimony, they may not be flexible enough to describe the complicated underlying relationship between the gene and the trait over time. Zhang and Zhong (2006) showed that when the parametric genetic effect function is misspecified, the power for detecting genetic effect can be greatly reduced to as low as 35%. It is therefore desirable to consider more flexible semiparametric models to analyze data from longitudinal genetic studies. To this end, Wu, Yang and Wu (2007) developed a nonparametric method using B-splines for the

QTL genotype-specific mean curve through a normal mixture model, and Zhao and Wu (2008) used a wavelet-based method. These mixture model based approaches may be difficult to implement in general extended pedigrees.

In this paper, we first present semiparametric regression model for overall polygenic effect with longitudinal data generated from family-based genetic studies such as the FHS. Next, we extend the polygenic model to accommodate genetic markers and model age-dependent associations. For family-based designs, subjects are nested in families and repeated measurements are nested within subjects. Therefore it is critical to account for both the subject-level and the measurement-level correlations in the statistical analyses to achieve more accurate estimation. One of the key features of the family-based longitudinal genetic studies is that subjects in the same family may not be independent, given any genetic marker. This is because the marker under consideration may not fully explain familial aggregation in a family. The residual unexplained genetic information aside from the marker is termed unspecified residual polygenic effect which is modeled as a random effect.

Mixed effects model naturally lends itself to account for residual polygenic effect between subjects in a family as well as serial correlation between repeated measurements of the outcome on the same subject. Meanwhile, it is desirable to model the baseline function and the genetic-effect function nonparametrically because there is usually limited information about the parametric forms of these functions. For this purpose, we propose to use penalized splines (P-spline; Eilers and Marx 1996) to estimate nonparametric functions in the model. Penalized splines based methods have become popular in the recent literature (Ruppert, Wand and Carroll 2003). In a penalized splines regression, an unknown smooth function is estimated by assuming a high-dimensional spline basis and imposing a penalty on the spline coefficients to control overfitting and achieve smooth fit. The number of knots in penalized splines regression is usually less than the sample size. Empirical and theoretical work has shown that the penalized spline as a reduced rank smoother can achieve similar quality of fit as full rank estimators such as smoothing splines (Ruppert 2002; Li and Ruppert 2008; Claeskens, Krivobokova, and Opsomer 2009). Another feature of penalized splines which makes it particularly suitable for analyzing longitudinal genetic data is its mixed model representation. By this representation, it is easy to handle random polygenic effect and all approaches developed here can be implemented by standard statistical software packages such as PROC MIXED in SAS or NLME in R, allowing researchers to use these methods on routine basis.

The primary interests in this work are to estimate baseline function, age-specific QTL effect and residual polygenic heritability, and to test for the QTL effect. The remaining of the paper is organized as follows. In section 2, we propose two semiparametric regression models for family-based longitudinal genetic studies to estimate baseline function and to test and estimate time-varying QTL effect and residual polygenic heritability. In section 3, we develop statistical methods for these two classes of models. In section 4, we perform simulation studies to investigate properties of the proposed methods. In section 5, we apply the developed methods to analyze the Framingham Heart Study blood pressure data. In section 6, we discuss implications of our findings on FHS and possible extensions of the proposed methods.

## 2 Models for longitudinal genetic studies

Here we propose semiparametric models for family-based longitudinal genetic studies. In section 2.1, we introduce partially linear mixed effects models to handle polygenic effect (overall genetic effect), and in section 2.2, we extend the models in section 2.1 to a semi-varying coefficient partially linear mixed effects model to incorporate the QTL effect.

## 2.1 Partially linear mixed effects model for polygenic effect

The first step in a genetic epidemiological study is to assess polygenic heritability of a trait by examining similarity of a trait in family members before using any genetic markers. The polygenic heritability quantifies the overall genetic impact on a trait. If there is no evidence of polygenic effect or familial aggregation, it may not be necessary to pursue further study such as linkage or association analysis that aims at locating the underlying loci affecting the trait. On the contrary, if the evidence for genetic factor predisposing a trait is observed, then in order to locate this factor along the genome, investigators decompose the polygenic effect into a major genetic effect at a specific locus and a residual polygenic effect contributed by other unlinked loci.

In statistical genetics theory, polygenic effect is treated as an unobserved random variable with covariance matrix specified by relationship between relatives (Lynch and Walsh 1998; Khoury, Beaty and Cohen 1993). To be specific, let $Y_{ijh}$ be the phenotype measurement for subject $j$ in family $i$ at visit $h$, and let $t_{ijh}$ denote the subject's age at this visit. A partially linear mixed effects model for $Y_{ij}(t_{ijh})$ is defined to be

$$Y_{ijh} = \mu(t_{ijh}) + x_{ij}^T \beta + \alpha_i + z_{ijh}^T \gamma_{ij} + \varepsilon_{ijh},$$
$$\alpha_i \sim N(0, \sigma_\alpha^2), \quad \gamma_i \sim N(0, \Sigma_i^\gamma), \quad \varepsilon_{ij} \sim GP(0, \vartheta_{ij}) \tag{1}$$

where $\mu(t)$ is an unspecified baseline function and $x_{ij}$ are environmental exposures such as sex with effects $\beta$, $\alpha_i$ are random shared environmental effects such as diet shared among family members, $\gamma_i = (\gamma_{i1}^T, \cdots, \gamma_{in_i}^T)^T$ are vectors of random polygenic effects, $z_{ijh}$ are design vectors for $\gamma_{ij}$ which can be time-dependent to capture an age-related polygenic effect, $\varepsilon_{ij} = (\varepsilon_{ij1}, \ldots, \varepsilon_{ijn_{ij}})^T$ are random measurement errors with possible serial correlation, and $GP(0, \vartheta_{ij})$ is a Gaussian process with covariance matrix $\vartheta_{ij}$. Inclusion of exposures with time-varying effects are deferred to the next section where we introduce time-varying QTL model. We assume that $\alpha_i$, $\gamma_{ij}$ and $\varepsilon_{ij}$ are independent. The random polygenic effect reflects overall genetic information in a trait. Their covariance structure depends on the relationship among family members (Khoury, Beaty and Cohen 1993, Chapter 7). Specifically,

$$\sum_i^\gamma = \text{Cov}(\gamma_i, \gamma_i^T) = 2K^i \otimes \Omega_\gamma, \quad \text{Cov}(\gamma_i, \gamma_{i'}^T) = 0 \text{ for } i \neq i', \tag{2}$$

where $K^i$ is a known kinship coefficient matrix whose $(j, j')$th element is determined by the relationship between subjects $j$ and $j'$ in family $i$, and $\Omega_\gamma$ is an unknown covariance of the polygenic effect. The kinship coefficient is defined as the probability of randomly drawing an allele in subject $j$ that is identical by descent (IBD) to an allele at the same locus randomly drawn from subject $j'$. For example, twice the kinship coefficient, $2K_{jj'}^i$, for a full sibling pair is 1/2 and for a half-sibling pair is 1/4 (Khoury, Beaty and Cohen 1993, page 211). Parameters in $\Omega_\gamma$ represent the unknown polygenic variance which we are interested in. The heritability is defined as the ratio of the genetic variance to the total variance, that is,

$$h_\gamma^2(t) = \sigma_\gamma^2(t) / \sigma_T^2(t), \tag{3}$$

where $\sigma_T^2(t) = \sigma_\alpha^2 + \sigma_\gamma^2(t) + \sigma_\varepsilon^2(t)$, $\sigma_\gamma^2(t) = \omega_{11} + 2\omega_{12}t + \omega_{22}t^2$ for a linear $z_{ijh}$ design vector, $\omega_{ij}$ is the $(i, j)$th element of $\Omega_\gamma$, and $\sigma_\varepsilon^2(t)$ is the variance of the residual random measurement

error. A test of $h_\gamma^2(t)$ based on functional principal components of heritability was proposed in Fang and Wang (2009).

Although the mixed effects model formulation of penalized splines allows the baseline and the QTL functions to be fitted by standard statistical softwares, one practical complication in genetic study is how to impose the correlation structure of polygenic effect predicted by kinship coefficients as shown in (2). In behavioral genetics, decomposing phenotypic variance into genetic and environmental components are typically done by structural equation models and estimated from specialized software such as Mx (Neale et al. 2004). Guo and Wang (2002) ignored the kinship correlation in order to use standard software to fit a multilevel model.

Rabe-Hesketh et al. (2008) showed that for most family designs, one can reparametrize the polygenic effect into a few family-specific and subject-specific random effects allowing for easy handling of polygenic effect by standard softwares. For example, for nuclear families, one replaces the polygenic effect $\gamma_{ij}$ in model (1) by two family-specific and a subject-specific random effects as

$$\gamma_{ij} = a_{i1}(M_{ij} + \frac{C_{ij}}{2}) + a_{i2}(F_{ij} + \frac{C_{ij}}{2}) + a_{ij}\frac{C_{ij}}{\sqrt{2}},$$

where $M_{ij}$ is a binary indicator for mother, $F_{ij}$ for father, and $C_{ij}$ for children. The family-specific random effects $a_{i1}$ and $a_{i2}$ induce required correlation between parents and each child and between the children. However the induced variance for children from these two random effects is only half of the desirable variance and the other half is induced by the subject-specific random effects $a_{ij}$. By this reparametrization, we can easily fit a semiparametric model with polygenic effect by a standard software.

## 2.2 Semi-varying coefficient partially linear mixed effects model for QTL effect

When genetic markers such as SNPs are available, we add marker genotypes to model (1) to assess association between a marker and a trait. Due to dense SNP genotyping, we assume that the QTL is either at the SNP marker under consideration or tightly linked to it.

Let $g_{ij}$ denote the SNP marker genotype for subject $j$ in family $i$ coded as the copies of minor alleles which takes value 0, 1, or 2. Let $x_{ij}$ denote time-invariant environmental covariate such as gender, and let $w_{ij}(t)$ denote time-varying exposures with potentially time-varying effect such as body mass index. A semi-varying coefficient partially linear mixed effects model for $Y_{ijh}$ is

$$Y_{ijh} = \mu(t_{ijh}) + x_{ij}^T\beta + \alpha_i + \tilde{\gamma}_{ij}^T z_{ijh} + \beta_g(t_{ijh})g_{ij} + \theta^T(t_{ijh})w_{ij}(t_{ijh}) + \varepsilon_{ijh}, \quad (4)$$

where $\tilde{\gamma}_{ij}$ is the residual polygenic effect aside from the QTL effect, and $\theta(t)$ is the coefficient vector for covariates $w_{ij}(t)$. In this model, in addition to the baseline function $\mu(t)$ and other covariate effects, we are interested in estimating the time-varying genetic function, $\beta_g(t)$.

The age-specific QTL heritability is then defined as (Falconer 1985)

$$h_g^2(t) = \sigma_g^2(t)/\sigma_T^2(t), \qquad (5)$$

where $\sigma_g^2(t) = \mathrm{Var}\{\beta_g(t)g_{ij}\} = \beta_g^2(t)\mathrm{Var}(g_{ij})$ and $\sigma_T^2(t) = \sigma_\alpha^2 + \sigma_\gamma^2(t) + \beta_g^2(t)\mathrm{Var}(g_{ij}) + \sigma_\varepsilon^2(t)$. The QTL heritability can be interpreted as the proportion of total variation explained by the QTL. The residual polygenic heritability contributed by other unlinked loci is $h_\gamma^2(t) = \sigma_\gamma^2(t)/\sigma_T^2(t)$. The total heritability in a trait is the sum of the QTL heritability and the residual polygenic heritability.

To test for association between a genetic marker and a trait, we consider the null hypothesis $H_0$: $\beta_g(t) = 0$. To test for constant genetic effect, that is, the genetic effect does not change over time, we consider the null hypothesis $H_0$: $\beta_g(t) = \beta_g$.

## 3 Statistical methods for longitudinal genetic studies

### 3.1 An estimation procedure for the partially linear mixed effects model

For simplicity, we use truncated polynomial basis in our estimation procedure. Extension to other basis such as B-splines is discussed in section 6. We approximate the mean function by a linear combination of spline basis functions

$$\mu(t;\eta) \approx \eta_0 + \eta_1 t + \cdots + \eta_q t^q + \sum_{m=1}^{M} \eta_{q+m}(t - \tau_m)_+^q,$$

where $\tau_m$, $m = 1, \cdots, M$ is a given sequence of knots and $q$ is the order of the splines. We discuss selection of knots later in this section. For given variance components, we estimate $\eta$ and $\beta$ by maximizing the penalized logarithm of the marginal likelihood defined as

$$-\frac{1}{2}\log|\sum| - \frac{1}{2}(Y - X\beta - W\eta)^T \sum{}^{-1}(Y - X\beta - W\eta) - \frac{1}{2}\lambda\eta^T J\eta, \qquad (6)$$

where $Y$ is a vector of outcome, $\Sigma$ is the covariance of $Y$, $J = \mathrm{diag}(\mathbf{0}_{q+1}, \mathbf{1}_M)$ is a penalty matrix, $X$ and $W$ are design matrices specified in the Appendix A.1 available at www.columbia.edu/~yw2016, and $\lambda$ is a smoothing parameter. When $\lambda$ goes to infinity, the spline coefficients are shrunk towards zero and the fit converges to a polynomial function. When $\lambda$ goes to zero, the fit converges to a weighted least square. The estimating equations for $\beta$ and $\eta$ are constructed in the Appendix A.1. The solution for $\eta$ takes the form of a ridge regression estimate.

It is well known that there is a mixed effects model representation of penalized splines (Ruppert, Wand and Carroll 2003; Wand 2003). We explore this connection to facilitate computation using standard software. For penalized splines, Wand (2003) showed that the solution to maximizing the penalized likelihood in (6) is identical to the best linear unbiased predictor (BLUP) from a linear mixed effects model with certain choice of smoothing parameter which we describe in the Appendix A.1. The key is to specify the spline coefficients $\eta_{q+1}, \cdots, \eta_{q+M}$ as random effects with the same variance and construct appropriate design matrices for the fixed and the random effects.

The tuning parameters for penalized splines include number and placement of knots and smoothing parameter $\lambda$. Once the number of knots has been chosen, we place them at equal sample quantiles of the observed $t_{ijh}$'s. The smoothness of the fit is controlled by both $M$

and $\lambda$. Ruppert (2002) and Claeskens, Krivobokova and Opsomer (2009) showed that when $M$ is adequately large, further increasing $M$ not only does not improve the fit but also can sometimes deteriorate the fit. For smooth and either monotonic or unimodal functions, moderate number of knots is usually sufficient (Ruppert 2002; Yu and Ruppert 2002). The smoothing parameter $\lambda$ controls overfitting for moderate to large number of knots and plays a more critical role than $M$.

For given $M$, $\lambda$ can be chosen by generalized cross validation (GCV), minimizing AIC or estimating by restricted maximum likelihood (REML). Krivobokova and Kauermann (2007) investigated behavior of several data-driven smoothing parameter selectors including REML and AIC with correlated data. It is found through theoretical derivation and simulations that when the correlation structure is misspecified, the AIC-based choice failed to estimate a function properly and the REML-based choice provides much more satisfactory fit and exhibits less variability (Krivobokova and Kauermann 2007). To accommodate possible misspecification of the correlation structure of $\varepsilon_{ij}(t)$, here we use REML to estimate the smoothing parameters as shown in the appendix.

### 3.2 An estimation procedure for the semi-varying coefficient linear mixed effects model

For model (4), we also approximate $\beta_g(t)$ by a linear combination of basis functions, that is,

$$\beta_g(t) \approx \xi_0 + \xi_1 t + \cdots + \xi_q t^q + \sum_{m=1}^{M} \xi_{q+m}(t - \tau_m)_+^q. \tag{7}$$

Varying-coefficients $\theta(t)$ for covariates other than genetic marker can be handled in a similar fashion by the approximation $\theta(t) \approx \theta_0 + \theta_1 t + \cdots + \theta_q t^q + \sum_{m=1}^{M} \theta_{q+m}(t - \tau_m)_+^q$. For given variance components, the penalized logarithm of the marginal likelihood of $\beta$, $\eta$ and $\xi$ is

$$-\frac{1}{2}\log|\sum| - \frac{1}{2}r'\sum{}^{-1}r - \frac{1}{2}\lambda_1 \eta' J\eta - \frac{1}{2}\lambda_2 \xi' J\xi - \frac{1}{2}\lambda_3 \theta' J\theta,$$

where $r = (Y - X\beta - W\eta - S_1\xi - S_2\theta)$, the design matrix $S_1$ and $S_2$ are defined in the Appendix A.2, and $\lambda_1$, $\lambda_2$, and $\lambda_3$ are smoothing parameters for the baseline, the genetic effect function and varying coefficient for other covariate, respectively. In the Appendix A.2, we expand the mixed effects model used to fit (1) to obtain the coefficients for time-varying genetic effect. As described there, we select $\lambda_1$, $\lambda_2$ and $\lambda_3$ by treating them as extra variance components and estimating by REML.

### 3.3 Estimating the total variance

Since the total phenotypic variance function $\sigma_T^2(t)$ is involved in the heritability function (3), a non-parametric estimation is desirable. Fan, Huang and Li (2007) proposed a semiparametric estimator of the covariance function $\vartheta(s, t)$. They assumed that the correlation function has a parametric form, that is, $\vartheta(s, t) = \text{Cov}(\varepsilon_{ij}(s), \varepsilon_{ij}(t)) = \rho_\varepsilon(s, t; \nu)$, where $\rho$ is a known function, and $\nu$ is a vector of parameters. They estimated the variance function $\vartheta(t, t) = \text{Var}(\varepsilon_{ij}(t))$ nonparametrically through local kernel smoothing. Here we propose a penalized splines based approach.

To be specific, we estimate the total variance function through a penalized splines regression based on the residuals from $Y_{ijh}$ after subtracting off the fitted mean curves, but not any of the variance components, therefore they retain the total variability in the outcome. Let let

$$\widehat{\varepsilon}_{ijh} = Y_{ijh} - \widehat{\eta}(t_{ijh}) - x_{ijh}\widehat{\beta} - \widehat{\beta}_g(t_{ijh})g_{ij} - \widehat{\theta}(t_{ijh})w_{ij}(t_{ijh}),$$

where $\hat{\eta}$, $\hat{\beta}_g$, and $\hat{\theta}$ are the fitted value of the mean and the QTL genetic function. Similar to the estimation of $\eta$ and $\beta_g$, we express $\log(\sigma_T^2(t))$ as a linear combination of basis functions,

$$\log(\sigma_T^2(t)) \approx \rho_0 + \rho_1 t + \cdots + \rho_q t^q + \sum_{m=1}^{M}\rho_{q+m}(t - \tau_m)_+^q.$$

We then estimate $\rho$ by fitting a penalized splines regression to $\log(\widehat{\varepsilon}_{ijh}^2)$. Using the estimated fixed coefficients and the BLUP of the random effects, the fitted value of the total variance function will be

$$\widehat{\sigma}_T^2(t_{ijh}) = \exp\{\widehat{\rho}_0 + \cdots + \widehat{\rho}_q t_{ijh}^q + \sum_{m=1}^{M}\widehat{\rho}_{q+m}(t_{ijh} - \tau_m)_+^q\}. \tag{8}$$

The estimated total variance is then used to calculate heritability in (3) and (5). We evaluate performance of this procedure through examining MASE, mean bias and variance of heritability estimates in section 4.

### 3.4 Testing for association between a marker and a trait

When the QTL genetic effect is time-invariant, the hypothesis of no association between a marker and a trait is specified by $H_0: \beta_g = 0$ versus $H_a: \beta_g \neq 0$, which can be examined by a standard Wald test. When fitting a time-varying QTL model, the hypothesis of no association is

$$H_0: \xi_0 = \xi_1 = \cdots \xi_q = 0, \text{ and } \sigma_\xi^2 = 0, \tag{9}$$

where $\xi_1, \cdots, \xi_q$ are coefficients for polynomial terms defined in (7) and $\sigma_\xi^2$ is the variance of the random spline coefficients $\xi_{q+1}, \cdots, \xi_{q+M}$ as described in the Appendix A.2. This hypothesis can be examined by a likelihood ratio test. Crainiceanu and Ruppert (2004) showed that the distribution of the likelihood ratio test of (9) for penalized splines mixed model is non-standard due to lack of independence and the variance component parameter is on the boundary under the null hypothesis. Using a conventional 50:50 mixture of chi-square distributions may be conservative. For mixed models involving one variance component, Crainiceanu and Ruppert (2004) derived the finite sample distribution and asymptotic distribution of the likelihood ratio test. However, the asymptotic distribution of likelihood ratio test for more complicated models with multilevel random effects is unknown. Greven et al. (2008) proposed to estimate the mixing proportion of chi-square distributions by simulation based on pseudo-likelihood ratio test for models with multiple variance components.

Here we propose to compute the *p* value by a permutation procedure. Since under the null hypothesis the marker genotypes are not associated to the trait, we can permute genotypes among subjects. However, it is not straightforward to randomize genotypes in a family sample because simple permutation would not maintain phenotype correlation among related individuals. The family members are not exchangeable under the null hypothesis because even there is no major QTL effect, there may exist residual polygenic effect causing family members to be correlated. Yang et al. (2010) proposed to permute genotypes among founders and then simulate offspring genotypes conditionally on permuted founders' genotypes based on Mendelian law while keeping the phenotypes as observed. Specifically, we first permute the genotypes of founders (subjects who do not have parents) in a pedigree. Give a set of permuted founders genotypes, we generate an offspring's genotype by randomly select an allele from each parent of the offspring following the Mendelian law. Genotypes of siblings in the same family are assigned independently given their permuted parental genotypes. For each copy of permuted genotype data, the same model fitting procedure is carried out. In a genome-wide association study (GWAS), it is computational challenging to conduct permutation for every SNP. Since the null distribution of the test statistic is the same for SNPs with the same founder genotype frequency with a given family data, one can group SNPs into strata that have the same or similar founder genotype frequency, and only one permutation null distribution is needed for each group (Yang et al. 2010).

## 4 Simulations

In this section, we investigate performance of our proposed estimation and testing procedures through Monte Carlo simulations. We simulated 100 nuclear families among which 50 had two children, 30 had three and 20 had four. The number of observations on each parent ranged from four to eight, the number of observations for children ranged from two to four, and each subject was examined every two or four years. The total number of observations was 1749. Subjects' age ranged from 10 to 75 with a mean of 39.5. These settings were close to the assessment schedule in the FHS. For the analysis involving genetic marker, we simulated a fully linked genetic polymorphism with a dominant effect and a minor allele frequency of 0.25. We assumed that the transmission of allele from parental generation to offspring generation follows Mendelian law.

### 4.1 Time invariant genetic effect

In the first few simulations, the baseline function $\mu(t)$ was a logarithm function, $-34.2 + 81.7 - \log(0.25 - (t + 21.7))$, where the parameters were estimated from fitting a logarithm function to the FHS cholesterol data. Such function was used to simulate the baseline and the genetic effect function on several traits at the Genetic Analysis Workshop 13 (GAW13, Daw et al. 2003), where the simulations were designed to mock the actual FHS data provided at the workshop. The random shared familial environmental factor $a_i$ had a variance of 16, and the polygenic effect $\gamma_{ij}$ had a variance of 4. These parameters were chosen so that the polygenic heritability is in the range of that estimated by the FHS investigators (Levy et al. 2000). The variance function of residuals was an exponential function, $\mathrm{Var}(\varepsilon_{ij}(t)) = \exp(0.02 - t)$. The correlation of the residuals was AR(1) with autocorrelation parameter 0.6. We also examine other functional forms of $\mu(t)$ such as the Gaussian or the sine function. The baseline function was estimated by cubic truncated polynomials with 15 knots.

In simulation setting 1, we assumed $\beta_g(t) = \beta_g$ in model (4), where the true values of $\beta_g$ are shown in table (2). We computed the mean average squared error (MASE) of $\hat{\mu}(t)$ as the mean across the 500 simulations of the average squared error,

$$\mathrm{ASE}(\mu) = \frac{1}{K} \sum_{t_k \in T_\kappa} [\widehat{\mu}(t_k) - \mu(t_k)]^2,$$

where $T_\kappa$ is a set of grid points over time and $K$ is the cardinality of $T_\kappa$. Define the MASE of $\widehat{h_\gamma^2}(t)$ and $\widehat{h_g^2}(t)$ similarly. We summarize the maximal absolute relative bias, the mean bias and the mean variance averaged over grid points $T_\kappa$, and the MASE in the fist five columns of Table 1 (setting 1), which showed a small relative bias and MASE. The estimated time-invariant genetic effect was 9.99 (true value: 10), with a mean estimated standard error of 0.26 (empirical standard error: 0.27).

We compare proposed semiparametric analyses where $\mu(t)$ was estimated through penalized splines with a correctly specified nonlinear mixed effects model analysis and a misspecified parametric analysis where $\mu(t)$ was assumed to be a quadratic polynomial function. The results were recorded in the last six columns of Table 1 (setting 1). As expected, it is evident that when $\mu(t)$ was misspecified its estimation had large bias. It may be of interest to note that misspecification of the baseline function also affects estimation of the heritabilities. The mean bias in the marker-specific and the total heritabilities ($\widehat{h_g^2}(t)$ and $\widehat{h_T^2}(t)$) increased by 43% and 54%, respectively, when $\mu(t)$ was misspecified. In terms of estimating the baseline function, the semiparametric estimators are less efficient than the parametric analysis under a correctly specified model. For the heritability estimators, the efficiency loss of the semiparametric estimators is less notable.

For the variance components, the estimated polygenic variance was 4.06 (true value: 4), and the family-specific variance component was 15.88 (true value: 16). The asymptotic distribution of the heritability estimates is not straightforward to derive due to definition of the heritability being the ratio of two non-independent variance estimators. To compute confidence interval, we use bootstrap resampling. As seen from Table 1, the maximal relative bias and MASE of the QTL heritability and the total heritability were small. We present the estimated marker-specific heritability, the total heritability and their confidence intervals in the left panel of Figure 1. The empirical and bootstrap standard errors were compared in the right panel of Figure 1. The bootstrap standard error tracked the empirical one closely.

Our next simulation experiments examine effects of different baseline function estimators on testing a genetic effect. We simulated data under the same model (4) with $\beta_g(t) = \beta_g$, various effect size of the genetic marker and various functional forms of $\mu(t)$ (see Table 2 for these specifications). The random measurement errors were simulated from a normal distribution with mean zero and variance 10. We tested the significance of $\hat{\beta}_g$ by a standard Wald test. We compare performance of the proposed semiparametric analysis where $\mu(t)$ is estimated through penalized splines with three other analyses: (1) Misspecifying $\mu(t)$ as a linear function; (2) Misspecifying $\mu(t)$ as a quadratic function; and (3) Correctly specifying $\mu(t)$ as a nonlinear function and estimating through fitting a non-linear mixed effects model. First, we examine the type I error of all four analyses. From the second, the sixth and the tenth row of Table 2, we see that the semiparametric analysis and the correctly specified nonlinear analysis maintains the nominal level of the type I error. However, the two misspecified analyses reported either substantially inflated or highly conservative type I error rate depending on the true form of $\mu(t)$ and how it is specified. For example, when the true baseline function is a sine function but misspecified as a linear or a quadratic polynomial, the type I error rate for a test for $\beta_g$ at 5% level can be as high as 0.99. The erroneous type I error may be explained by two reasons: First, incorrect estimation of the baseline function

under a misspecified model may lead to incorrect standard error estimate of $\hat{\beta}_g$; Second, the mean of $\mu(T)$ across observed time points is not a constant across different genotype groups, i.e., $E(\mu(T)|G = g)$ differs across levels of $G$ in a partially linear model which may lead to inconsistent estimate of $\hat{\beta}_g$.

Next we compare power of the test when $\mu(t)$ is estimated nonparametrically with the correctly specified nonlinear analysis. From Table 2, we see that the power for testing genetic effect is slightly larger with a correctly specified nonlinear baseline function comparing to the semiparametric analysis, with a difference up to 5%. For the scenarios in Table 2 where the two misspecified analyses had conservative type I error, we also examined their power. As expected, the power was greatly reduced with a power loss up to 95% comparing to the semiparametric analysis. For example, when the true $\mu(t)$ is a Gaussian function but misspecified as a linear or a quadratic function, the power for detecting a genetic effect was zero. In addition to a highly conservative type I error rate, this may also be due to substantial increase of variability of the estimator $\hat{\beta}_g$ when the baseline function was misspecified in these cases.

To summarize, the first set of the simulations suggest that misspecification of the baseline function has a non-ignorable effect on the type I error of testing the genetic effect even when the genetic effect does not change with time. Furthermore, the power of testing $\beta_g$ when treating $\mu(t)$ as a nonparametric function is comparable to correctly specifying $\mu(t)$ as a nonlinear function.

## 4.2 Time varying genetic effect

The second simulation setting examines properties of our methods when $\beta_g(t)$ changes with time. The performance of the baseline function estimator was comparable to the time-invariant case (Table 1, setting 2). From Table 1 (setting 2), we see that the time-varying genetic effect $\hat{\beta}_g(t)$ was estimated well with small MASE. We show the true and the estimated genetic effect and its confidence interval in the left panel of Figure 2. The bootstrap standard error and the empirical standard error shown in the right panel of Figure 2 were very close. The age-specific QTL heritability and total heritability were estimated well with the maximal relative bias 0.02 and 0.01, respectively (Table 1, setting 2). The bootstrap and empirical standard error were close which suggests a satisfactory performance of the bootstrap procedure on assessing variabilities of heritability estimates (the corresponding figure is similar to Figure 1 and not shown).

Similar to the previous section, we compare the estimation bias and MASE of $\hat{\beta}_g(t)$ in a semiparametric analysis with a misspecified parametric analyses where we assumed $\beta_g(t)$ to be a quadratic polynomial and with a correctly specified nonlinear mixed effects model analysis. In all analyses, we kept the estimation of $\mu(t)$ nonparametric because the analyses in the previous section showed a profound effect of misspecifying $\mu(t)$ on testing $\beta_g$. From Table 1 (setting 2), we see that the mean bias of $\hat{\beta}_g(t)$ over time with a misspecified quadratic model increased from 0.036 in a nonparametric method to 0.23. The mean bias of the estimated marker specific heritability, $\widetilde{h_m^2}(t)$, increased from 0.004 to 0.34, which is substantial. The mean bias of the estimated total heritability increased from 0.006 to 0.37.

The rest of the simulations concern testing of $\beta_g(t)$. The random measurement errors were simulated from a normal distribution with mean zero and variance six. The hypothesis $\beta_g(t) = 0$ was tested by the permutation procedure described in Section 3.5 in the semiparametric analysis. In all analyses, the baseline function was again estimated nonparametrically. We examine several functional forms for $\beta_g(t)$ including logarithm, Gaussian and sine. The Gaussian function was used to model genetic effect on blood pressure in Shi and Rao (2008).

We first examine the type I error of the semiparametric analysis and two parametric analyses under a misspecified model. From Table 3, we see that all three analyses maintain the correct nominal level of type I error. We then examine power of testing $\beta_g(t) = 0$. Again as expected, we see from table 4 that the power is greatest for the nonlinear mixed effects model analysis with a correctly specified model. However, in the real applications such a true function is unknown and the computational algorithm in a nonlinear analysis may not converge in many cases especially when starting values are poor or the sample size is small or moderate. Comparing the semiparametric approach to the misspecified parametric approaches, the power loss for the latter ranges from 0% to 55%. The power loss was more substantial for the Gaussian and the sine function comparing to the logarithm function which suggests that the power depends on the unknown functional form of the true genetic effect and the assumed parametric model. For the genetic effect that changes with time but has an average effect of zero across all time points (i.e., $\frac{1}{\sum_{ij} T_{ij}} \sum_{ijh} \beta_g(t_{ijh}) \approx 0$), the linear or quadratic analysis has very low power (close to zero) to detect the genetic effect.

To summarize, these simulations suggest that misspecifying $\beta_g(t)$ while estimating $\mu(t)$ nonparametrically does not affect type I error rate of testing $\beta_g(t) = 0$, but may reduce power substantially.

## 5 Application to the Framingham Heart Study

In this section, we apply the proposed methods to analyze the Framingham Heart Study longitudinal blood pressure (BP) data and SNP genotype data. High blood pressure (BP) is considered as a major risk factor for stroke and heart disease and it affects about one-third of the US adult population (Levy et al. 2009). Systolic and diastolic blood pressure (SBP and DBP) are complex traits that may be influenced by both environmental and genetic factors. The heritability of systolic blood pressure is estimated to be high (30% to 60%, Levy et al. 2000), which suggests a substantial genetic contribution. Recently, large-scale genome-wide association studies (GWAS) have emerged as powerful tools to identify genes associated with complex traits such as BP. Levy et al. (2009) performed a prospective meta-analysis on six GWAS including the FHS and identified multiple SNPs significantly associated with SBP and DBP at the genome-wide significance level. However, nonparametric estimation of age-specific QTL effect or time-varying polygenic of BP has not been examined in the literature. We analyze a subset of the FHS subjects (about 6000 subjects) and a subset of SNPs in four candidate regions.

In the FHS, the phenotype and the genotype data are collected from three cohorts. The Original Framingham Cohort (Cohort 1) was first examined in 1948 and has been examined every two years thereafter. The Offspring Cohort (Cohort 2), composed primarily of offspring of the original cohort and the spouses of these offspring, was examined first in 1971 and has been examined approximately every four years using protocols similar to those used for study of the Original Cohort. Between 2002 and 2005 the study enrolled the Third generation (Gen3) of the Framingham Heart Study. At each exam, the physician measured systolic and diastolic blood pressure twice and the average of the two measurements was used as the phenotype in the analysis.

Although the FHS started at an era when no antihypertensive treatment was available, as the study progressed, antihypertensive treatment became available and was prescribed to some of the subjects with hypertension. It is known that the treatment effect is a confounder for genetic effect which may lead to underestimated genetic effect without any adjustment (Levy et al. 2000, Tobin et al. 2005). Tobin et al. (2005) examined bias and variance of ten methods on adjusting for treatment effect and found that one of the best methods is to add a

reasonable number to observed SBP for subjects on antihypertensive treatment. Following Tobin et al. (2005) and Levy et al. (2009), we added 10 mm Hg to observed SBP values and 5 mm Hg to observed DBP values for participants who were taking treatment.

We restricted our analysis to observations between age 30 and 75. The total sample size in our analysis was 6082 from 930 pedigrees (including 2934 nuclear families) and the mean number of subjects was 6.54 per extended family. There were 14505 records and each subject had an average of 2.38 measurements of SBP and DBP, respectively. The age of the participants at the first visit ranged from 25 to 72. The mean age for all subjects at all visits was 45.7 years. The mean observed SBP was 121.2 mm Hg and the mean observed DBP was 76.1 mm Hg. There were 11% subjects on antihypertensive treatment in at least one exam and 12% of observations were taken when subjects were on treatment. The mean body mass index (BMI) was 23.54.

In all our analyses, we included gender as a covariate with time-invariant effect and BMI as a time-varying covariate with varying coefficient. We estimated the baseline function by a cubic truncated polynomial with ten knots. We split pedigrees into nuclear families for easy handling of familial correlations. We first computed the baseline function and the polygenic heritability without using SNP markers as in model (1). The estimated age-specific baseline function and its 95% confidence interval are superimposed on SBP measurements of 300 randomly selected subjects in the left panel of Figure 3. There is an increasing trend of mean SBP over time. The mean SBP was 123.5 mm Hg (CI: 122.9, 124.2) at age 30 and increased to 138.6 mm Hg (CI: 134.8, 142. 4) at age 75. The corresponding plot for DBP was shown in the right panel of Figure 3. The polygenic heritability of SBP and its confidence interval are shown in the left panel of Figure 4. Heritability was highest at age 35.4 and it then decreased to 0.44 (CI: 0.40, 0.50) at age 50 and 0.23 at age 65 (CI: 0.18, 0.27). The long term average heritability was reported to be between 0.3 and 0.6 (Levy et al. 2000, Levy et al. 2009), which is in the range of our age-specific estimates. The total variance function increases over time and is presented in the right panel of Figure 4. For DBP, the polygenic heritability decreases with age. It was 0.44 (CI: 0.37, 0.54) at age 35.4 and then decreased to 0.29 (CI: 0.17, 0.47) at age 75. Theses heritabilities and the variance function of DBP are presented in Figure 5. Overall, DBP exhibits lower heritability than SBP. The gender effect was estimated as 1.57 (CI: 0.80, 2.34) with men having higher SBP, on average.

Levy et al. (2009) conducted meta analysis of six GWAS of blood pressure and reported several promising regions which may harbor genes predisposing BP. We selected four promising candidate regions containing significant SNPs reported in Levy et al. (2009) to analyze. There were 265 SNPs in the four regions among which 109 SNPs were from two regions on chromosome 12 (86 from region 88300Kb to 88800Kb and 23 from region 110200Kb to 110600Kb), 104 were from a region on chromosome 11 (16600Kb to 17100Kb) and 52 were from a region on chromosome 3 (41700Kb to 42100Kb). Each of these regions spans about 500Kb on a chromosome. We first fit a time-invariant model with a nonparametric baseline function but a constant genetic effect, that is, $\beta_g(t) = \beta_g$ in model (4). Since adjusting for multiple comparisons by Bonferroni correction is conservative for dense SNPs in linkage disequilibrium (LD), we use methods proposed in Gao et al. (2008). Specifically, we use principal components analysis to compute effective number of SNPs needed to explain 99.5% of variability of all 234 SNPs and then divide the overall significance level (0.05) by this number. The resulting effective number of SNPs needed is 104, and the adjusted significance level is $4.81 \times 10^{-4}$. There was one SNP on chromosome 12 significant for SBP at this level and none for DBP (Table 5).

In addition to the time-invariant analysis, we also fit a time-varying genetic effect model and test for the hypothesis (9) on all SNPs. We found four significant SNPs for DBP and five for

SBP after adjusting for multiple comparisons. None of these SNPs were identified through the time-invariant analyses. For some SNPs, their $p$ values in the time-invariant model showed suggestive results for association (for example, the $p$ value for rs1052501 in a time-invariant analysis was 0.002), however, they did not reach the significance level. Other SNPs would not have been identified from a time-invariant analysis (for example, the $p$ value for rs10858911 in a time-invariant analysis was 0.32). As an example, we show the age-specific effects and their confidence intervals of two SNPs in Figure 6. The SNP rs1052501 is in LD with three other SNPs in the same region identified through the time-varying analysis for the DBP. Two SNPs identified for SBP, rs4757448 and rs17700056, are in LD. The time-varying analysis suggests that there may be genes not only affect the long term average SBP, but also affect change of SBP with time. We discuss implications of these findings and compare with the time-invariant analysis in the next section.

## 6 Discussion

In this work, we propose semiparametric regression analysis of genetic studies with longitudinal phenotypes by penalized splines. Although the age-specific QTL effect function is modeled non-parametrically, the test of no association is examined by only a few parameters in hypothesis (9). Mixed effects model representation of penalized splines provides a convenient way to handle polygenic effect and shared environmental effect in genetic studies. Our simulations show that misspecifying the baseline function in a parametric analysis has a substantial effect on type I error rate and power of testing the genetic effect regardless of whether it changes with time. Furthermore, when the true genetic effect is a constant, the semiparametric analysis has comparable power comparing to a nonlinear analysis under a correctly specified model of the baseline function. It is therefore beneficial to model the baseline function nonparametrically. Misspecifying the genetic effect when the true effect varies with time in a parametric analysis can reduce power significantly, especially when the average genetic effect over time is small. The proposed semiparametric procedure provide an alternative to existing dominating time-invariant analysis and parametric linear or quadratic model for longitudinal genetic designs.

Although here the statistical procedures are developed for longitudinal data, they are also applicable to cross-sectional data when subjects' age is recorded. Let $t_{ij}$ denote the age of the $i$th subject in the $j$th family. A model similar to (4) for cross-sectional data is In addition, for population based case-control studies, the outcome is a binary variable. Penalized splines regression introduced here can be extended to generalized outcome through a connection with generalized mixed effects model as discussed in Ruppert, Wand and Carroll (2003).

Population stratification is a potential confounder in genetic association studies. However, for FHS all the study subjects are recruited from Framingham, Massachusetts where the majority of the population is Caucasian and population stratification is found to be negligible (Wilks et al. 2005). Nevertheless, one approach to adjust for population admixture is to estimate it by a principal components analysis and include the first few principal components as covariates in the model (Price et al. 2006), which can be readily incorporated in the framework of our proposed methods. The principal component weights are computed from founders in families and projected onto offsprings to create principal components scores which are then included in a regression analysis. Another method is to incorporate the permutation procedure implemented in the FBAT (family-based association test, Rabinowitz and Laird 2000) to the our permutation test of the genetic effect. To be specific, one permutes offspring's genotypes given minimal sufficient statistics of the genetic model under the null. A third strategy to adjust for population admixture in a regression based analysis with family data is to include expected value of the genotype-related covariates given the minimal sufficient statistic for the genetic model under the null as additional

covariates (Yang et al. 2000). This approach is the estimation analogy of the FBAT. Wang et al. (2010) discussed an improvement of Yang et al. (2000) that computes the optimal covariate to minimize the estimation variance and include these covariates in a regression analysis.

We implemented our methods with truncated polynomial basis. Other basis such as B-splines can also be used. Models based on B-splines are equivalent to truncated polynomials through a re-parametrization. The penalty matrix in (6) for B-splines, however, does not have the simple ridge penalty form and needs to be adapted. Eilers and Marx (1996) proposed a difference-based penalty. Wand and Ormerod (2008) considered a penalty matrix that is a direct generalization of smoothing splines (O'Sullivan penalized splines) and provided mixed model representation. These works allow our methods to be extended to B-splines.

Although the proposed methods are illustrated through a candidate region analysis of the FHS data, the mixed effects model based semiparametric analysis can be implemented for analyses at a much larger scale, for instance, a GWAS. In our application of the FHS data with 6082 subjects and 14505 observations, on average for each SNP the proposed procedure took 1.5 minutes to run on a Dell desktop with 2.00GHz CPU and 3.25GB memory using R package "lme". To complete a GWAS with 500K SNPs, this amounts to 2.6 days on a computing cluster with 200 nodes each with a 2.00GHz CPU or about 10 days for a cluster of 50 nodes. In our experience, SAS procedure "mixed" appears to improve computational efficiency in some cases up to 20%.

Our analyses identified six SNPs for SBP and four SNPs for DBP residing in three genes. The SNP rs11065951 locates within the gene ATXN2, which is a cytoplasmic protein. Lastres-Becker (2008) found that ATXN2 knock-out mice exhibited reduced fertility, locomotor hyper-activity, and abdominal obesity and hepatosteatosis at the age of 6 months. ATXN2 was also reported to associate with neurological disorders (Huynh et al. 1999), renal functions (Kottgen et ak. 2010) and obesity (Figueroa KP et al. 2009) which may share some pathway with BP. Four SNPs (rs4757448, rs17700056, rs7943587, and rs7121911) locate in a protein coding gene, PLEKHA7, which was reported to be linked to blood pressure at genome-wide level in another joint meta-analysis of GWAS studies for blood pressure (CHARGE and Global BPgen, Newton-Cheh et al. 2009). Four linked SNPs (rs1052501, rs7648578, rs2128834 and rs3774372) were located in the gene ULK4, which is an Unc-51-like kinase. This gene was also identified in CHARGE study as a candidate locus for blood pressure. However, little is reported on the relationship between function of this gene and blood pressure. Gene expression analysis has confirmed SNPs in ULK4 alter gene expression levels in liver and lymphoblastoid cell lines (Levy et al. 2009). Our analysis showed a potential time-varying effect at this locus which may deserve further functional research.

Aging is a complex process during which many biological and physiological changes take place which in turn may change a range of phenotypes, including blood pressure, and may change the interplay between the environmental and genetic factors. Therefore age may represent a surrogate of constellations of unmeasured factors. Taking into account of the gene-age interaction in a genetic association study may help overcome some of inconsistencies in replicating a genetic finding and boost power (Lasky-Su et al. 2008). Our time-invariant analysis identifies two SNPs for SBP and the time-varying analysis identifies a distinct SNP for SBP and four SNPs for DBP. None of the SNPs was identified by both analyses. Some of the SNPs may be missed if only the time-invariant analysis was carried out. These results illustrate the complementary feature of the two analyses. When the true genetic effect does not vary with time, a time-invariant model may identify more SNPs due

to parsimony of the model. However, when the genetic effect does change with time or when age acts as a surrogate of unmeasured factors causing varying genetic effect, failure to acknowledge the time trend may reduce power or lead to irreproducible results (Shi and Rao 2008, Lasky-Su et al. 2008).

Despite large efforts on gene mapping through GWAS, until recently there were few known genetic variants found to be reproducibly associated with common disease. Part of the inconsistency may be explained by the dominant time-invariant analyses strategy (Laksy-Su et al. 2008). The general semiparametric approaches we develop here may be applied to model age-dependent genetic effects leading to more powerful genetic data analysis and potentially more consistent results. Our analyses results also suggest new hypothesis of possible time-varying genetic effect on blood pressure at several loci which needs to be confirmed by future larger study. In addition, estimating age-specific heritability and genetic-effect has implications for designing subsequent studies and developing treatment of a disease: sampling subjects at the age where heritability is at its peak would enhance power of an association study, which is very important for detecting genes with moderate effects; designing a future GWAS rests on accurate estimation of potential time-varying effect size of a gene; and interventions may target different environmental or genetic factor at different age depending on which factor is dominant.
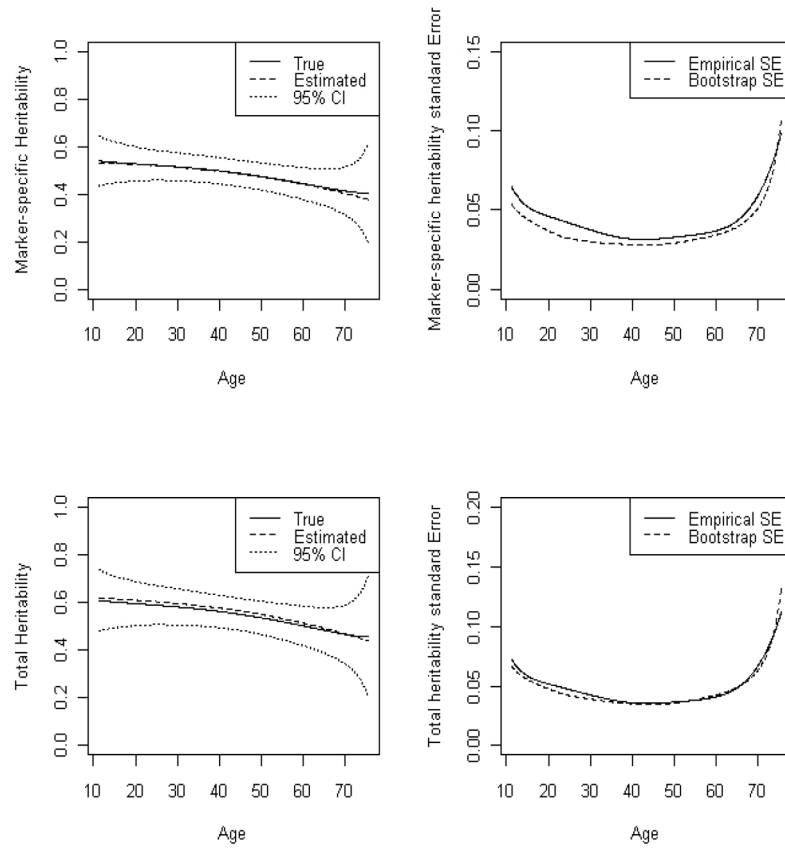
## Acknowledgments

## References

Atwood L, Heard-Costa N, Cupples L, Jaquish C, Wilson P, D'Agostine R. Genome-wide Linkage Analysis of Body Mass Index Across 28 Years of the Framingham Heart Study. American Journal of Human Genetic. 2002; 71:1044–1050.

Claeskens G, Krivobokova T, Opsomer J. Asymptotic Properties of Penalized Spline Estimators. Biometrika. 2009; 96:529–544.

Crainiceanu C, Ruppert D. Restricted Likelihood Ratio Tests in Nonparametric Longitudinal Models. Statistica Sinica. 2004; 14:713–729.

Daw EW, Morrison J, Zhou X, Thomas D. Genetic Analysis Workshop 13: Simulated Longitudinal Data on Families for a System of Oligogenic Traits. BMC Genetics. 2003; 4(Suppl 1):S3. [PubMed: 14975071]

Dawber TR, Meadors GF, Moore FEJ. Epidemiological Approaches to Heart Disease: the Framingham Study. American Journal of Public Health. 1951; 41:279–286. [PubMed: 14819398]

Eilers P, Marx B. Flexible smoothing with B-splines. Statistical Science. 1996; 11:89–121.

Fan J, Huang T, Li R. Analysis of Longitudinal Data with Semiparametric Estimation of Covariance Function. Journal of the American Statistical Association. 2007; 102:632–641. [PubMed: 19707537]

Falconer, DS. Introduction to Quantitative Genetics. 2. New York: Longman; 1985.

Fang Y, Wang Y. Testing for genetic effect on functional traits by functional principal components analysis based on heritability. Stat Med. 2009; 28(29):3611–3625. [PubMed: 19731232]

Figueroa KP, Farooqi S, Harrup K, Frank J, O'Rahilly S, Pulst SM. Genetic variance in the spinocerebellar ataxia type 2 (ATXN2) gene in children with severe early onset obesity. PLoS One. 2009; 4(12):e8280. [PubMed: 20016785]
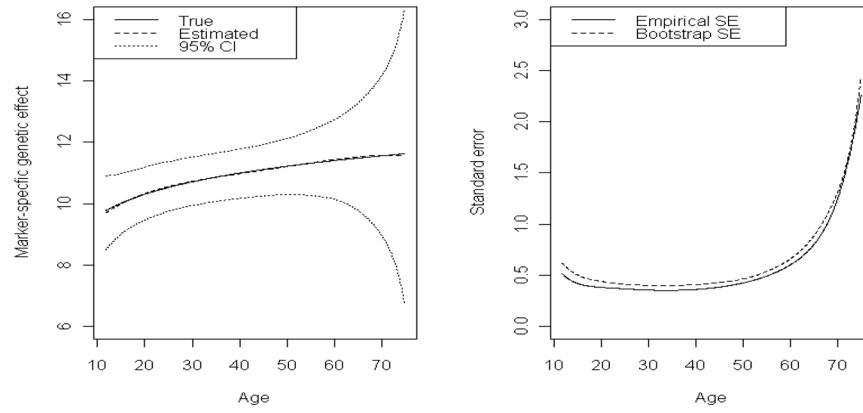
Gao X, Starmer J, Martin E. A Multiple Testing Correction Method for Genetic Association Studies Using Correlated Single Nucleotide Polymorphisms. Genet Epidemiol. 2008; 32:361–369. [PubMed: 18271029]

Greven S, Crainiceanu C, Kchenhoff H, Peters A. Restricted Likelihood Ratio Testing for Zero Variance Components in Linear Mixed Models. Journal of Computational and Graphical Statistics. 2008; 17(4):870–891.

Gudbjartsson DF, Bjornsdottir US, Halapi E, Helgadottir A, Sulem P, Jonsdottir GM, Thorleifsson G, Helgadottir H, et al. Sequence Variants Affecting Eosinophil Numbers Associate with Asthma and Myocardial Infarction. Nature Genetics. 2009; 41:342–347. [PubMed: 19198610]

Guo G, Wang J. The Mixed Model or Multilevel Model for Behavior Genetics Analysis. Behavior Genetics. 2002; 32:37–49. [PubMed: 11958541]

He Q, Berg A, Li Y, Vallejos CE, Wu R. Mapping Genes for Plant Structure, Development and Evolution: Functional Mapping Meets Ontology. Trends Genet. 2010; 26:39–46. [PubMed: 19945189]

Huynh DP, Del Bigio MR, Ho DH, Pulst SM. Expression of ataxin-2 in brains from normal individuals and patients with Alzheimer's disease and spinocerebellar ataxia 2. Annals of Neurology. 1999; 45(2):232–41. [PubMed: 9989626]

Jarvik GP, Goode EL, Austin MA, Auwerx J, Deeb S, Schellenberg GD, Reed T. Evidence that the Apolipoprotein E-Genotype Effects on Lipid Levels Can Change with Age in Males: a Longitudinal Analysis. American Journal of Human Genetics. 1997; 61:171–181. [PubMed: 9245998]

Kauermann, G.; Wegener, M. Functional Variance Estimation using penalized splines with principal components analysis. Statistics and Computing. 2009. http://www.springerlink.com/content/09603174/

Khoury, M.; Beaty, H.; Cohen, B. Fundamentals of Genetic Epidemiology. New York: Oxford University Press; 1993.

Kttgen A, Pattaro C, Bger CA, Fuchsberger C, Olden M, Glazer NL, Parsa A, Gao X, Yang Q, et al. New loci associated with kidney function and chronic kidney disease. Nat Genet. 2010; 42(5): 376–384. [PubMed: 20383146]

Krivobokova T, Kauermann G. A Note on Penalized Splines with Correlated Errors. Journal of the American Statistical Association. 2007; 102(480):1328–1337.

Lasky-Su J, Lyon HN, Emilsson V, Heid IM, Molony C, Raby BA, Lazarus R, Klanderman B, Soto-Quiros ME, Avila L, Silverman EK, et al. On the Replication of Genetic Associations: Timing Can Be Everything! American Journal of Human Genetics. 2008; 82:849–858. [PubMed: 18387595]

Lastres-Becker I, Brodesser S, Ltjohann D, Azizov M, Buchmann J, Hintermann E, Sandhoff K, Schrmann A, Nowock J, Auburger G. Insulin Receptor and Lipid Metabolism Pathology in Ataxin-2 Knock-out Mice. Human Molecular Genetics. 2008; 17(10):1465–1481. [PubMed: 18250099]

Levy D, DeStefano AL, Larson MG, O'Donnell CJ, Lifton RP, Gavras H, Cupples LA, Myers RH. Evidence for a Gene Influencing Blood Pressure on Chromosome 17. Genome Scan Linkage Results for Longitudinal Blood Pressure Phenotypes in Subjects from the Framingham Heart Study. Hyperension. 2000; 36:477–483.

Levy D, Ehret GB, Rice K, Verwoert GC, Launer LJ, Dehghan A, Glazer NL, Morrison, et al. Genome-wide Association Study of Blood Pressure and Hypertension. Nature Genetics. 2009; 41:677–687. [PubMed: 19430479]

Li Y, Ruppert D. On the Asymptotics of Penalized Splines. Biometrika. 2008; 95:415–436.

Lynch, M.; Walsh, B. Genetics and Analysis of Quantitative Traits. Massachusetts: Sinauer Associates; 1998.

Meng W, Mushika Y, Ichii T, Takeichi M. Anchorage of Microtubule Minus Ends to Adherens Junctions Regulates Epithelial Cell-cell Contacts. Cell. 2008; 135:948–959. [PubMed: 19041755]

Neale, MC.; Maes, HH. Mx: Statistical Modeling. 6. Rchomond, VA: Virginia Common Wealth University, Department of Psychiatry; 2004. Downloadable from: http://www.vcu.edu/mx/

Newton-Cheh C, Johnson T, Gateva V, Tobin MD, Bochud M, Coin L, Najjar SS, Zhao JH, et al. Genome-wide association study identifies eight loci associated with blood pressure. Nature Genetics. 2009; 41(6):666–676. [PubMed: 19430483]

Price AL, Patterson NJ, Plenge RM, Weinblatt ME, Shadick NA, Reich D. Principal Components Analysis Corrects for Stratification in Genome-wide Association Studies. Nature Genetics. 2006; 38(8):904–9. [PubMed: 16862161]

Province MA, Rao DC. Path Analysis of Family Resemblance with Temporal Trends: Applications to Height, Weight, and Quetelet index in Northestern Brazil. American Journal of Human Genetics. 1985; 37:178192.

Rabinowtiz D. Adjusting for Population Heterogeneity and Misspecified Haplotype Frequencies When Testing Nonparametric Null Hypotheses in Statistical Genetics. Journal of the American Statistical Association. 2002; 92:742–758.

Rabinowitz D, Laird N. A unified approach to adjusting association tests for population admixture with arbitrary pedigree structure and arbitrary missing marker information. Human Heredity. 2000; 50:211–223. [PubMed: 10782012]

Rabe-Hesketh S, Skrondal A, Gjessing HK. Biometrical modelling of twin and family data using standard mixed model software. Biometrics. 2008; 64:280–288. [PubMed: 17484777]

Raff RA. Evo-devo: the evolution of a new discipline. Nat Rev Genet. 2000; 1:74–79. [PubMed: 11262880]

Rice SH. A general population genetic theory for the evolution of developmental interactions. Proc Natl Acad Sci. 2002; 99:15518–15523. [PubMed: 12438697]

Ruppert D. Selecting the Number of Knots for Penalized Splines. Journal of Computational and Graphical Statistics. 2002; 11:735–757.

Ruppert, D.; Wand, MP.; Carroll, RJ. Semiparametric Regression. New York: Cambridge University Press; 2003.

Shi G, Rao DC. Ignoring Temporal Trends in Genetic Effects Substantially Reduces Power of Quantitative Trait Linkage Analysis. Genetic Epidemiology. 2008; 32:61–72. [PubMed: 17703462]

Strauch J, Golla A, Wilcox MA, Baur MP. Genetic Analysis of Phenotypes Derived from Longitudinal Data: Presentation Group 1 of Genetic Analysis Workshop 13. Genetic Epidemiology. 2003; 25(Suppl 1):S5–S17. [PubMed: 14635164]

Tobin MD, Sheehan NA, Scurrah KJ, Burton PR. Adjusting for Treatment Effects in Studies: Antihypertensive Therapy and Systolic Blood Pressure. Statistics in Medicine. 2005; 24:2911–2935. [PubMed: 16152135]

Van Steen K, McQueen MB, Herbert A, Raby B, Lyon H, Demeo DL, Murphy A, Su J, Datta S, Rosenow C, et al. Genomic screening and replication using the same data set in family-based association testing. Nat Genet. 2005; 37:683–691. [PubMed: 15937480]

Wand MP. Smoothing and Mixed Models. Computational Statistics. 2003; 18:223–249.

Wand MP, Ormerod JT. On O'Sullivan Penalised Splines and Semiparametric Regression. Australia and New Zealand Journal of Statistics. 2008; 50:179198.

Wang Y, Yang Q, Rabinowitz D. Unbiased and Efficient Estimation of the Effect of Candidate Genes on Quantitative Traits in the Presence of Population Admixture. Biometrics. 2010 In press.

Wilk JB, Manning AK, Dupuis J, Cupples LA, Larson MG, Newton-Cheh C, Demissie S, DeStefano AL, Hwang SJ, Liu C, Yang Q, Lunetta KL. No Evidence of Major Population Substructure in the Framingham Heart Study. Genetic Epidemiology. 2005; 29:286.

Wu, R.; Ma, C.; Casella, G. Statistical Genetics of Quantitative Traits: Linkage, Maps, and QTL. New York: Springer; 2007.

Wu S, Yang J, Wu R. Semiparametric Functional Mapping of Quantitative Trait Loci Governing Long-Term HIV Dynamics. Bioinformatics. 2007; 23:i569–i576. [PubMed: 17646344]

Yang Q, Rabinowitz D, Isasi C, Shea S. Adjusting for confounding due to population admixture when estimating the effect of candidate genes on quantitative traits. Human Heredity. 2000; 50:227–233. [PubMed: 10782014]
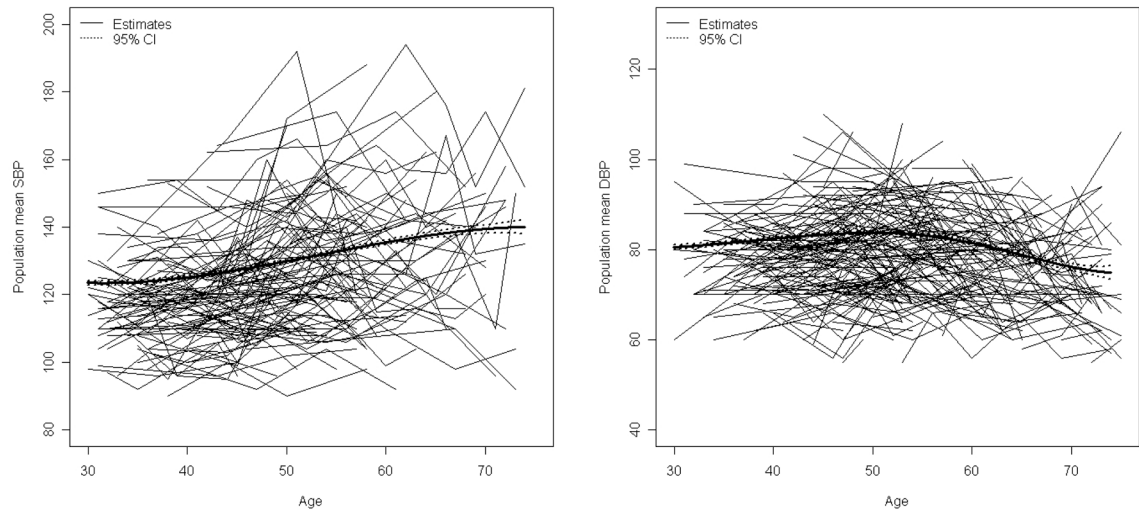
Yang Q, Wu H, Guo C, Fox C. Analyze Multivariate Phenotypes in Genetic Association Studies by Combining Univariate Association Tests. Genet Epidemiol. 2010; 34(5):444–454. [PubMed: 20583287]

Wang Y, Yang Q, Rabinowitz D. Unbiased and efficient estimation of the effect of candidate genes on quantitative traits in the presence of population admixture. Biometrics. 2010 In press.

Zhang H, Zhong X. Linkage Analysis of Longitudinal Data and Design Consideration. BMC Genetics. 2006; 7:37. [PubMed: 16768806]

Zhao W, Zhu J, Gallo-Meagher M, Wu R. A Unified Statistical Model for Functional Mapping of Environment-Dependent Genetic Expression and Genotype × Environment Interactions for Ontogenetic Development. Genetics. 2004; 168:17511762.

Zhao W, Wu R. Wavelet-Based Nonparametric Functional Mapping of Longitudinal Curves. Journal of American Statistical Association. 2008; 103:714–725.
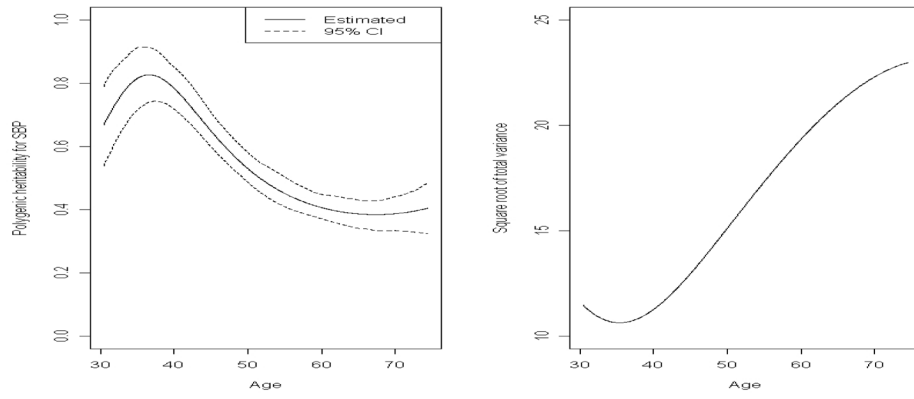
**Figure 1.**
Time-invariant genetic effect model: Estimated marker-specific heritability (top panel) and total heritability (bottom panel)
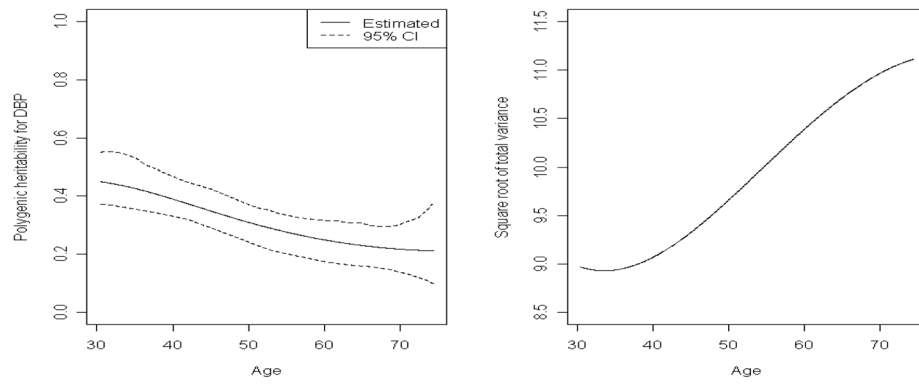
**Figure 2.**
Time-varying genetic effect model: Estimated age-specific genetic effect (left panel) and the bootstrap and the empirical standard error (right panel)

**Figure 3.**
SBP (left) and DBP (right) over time for 300 randomly selected subjects in the FHS and the estimated population mean function
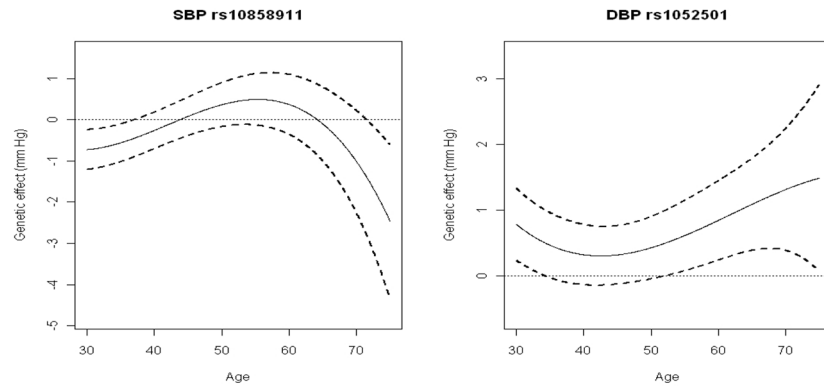
**Figure 4.**
Age-specific polygenic heritability of SBP (left panel) and total variance function (right panel)

**Figure 5.**
Age-specific polygenic heritability of DBP (left panel) and total variance function (right panel)

**Figure 6.**
Age-specific effects of two significant SNPs identified from the time-varying analysis

**Table 1**

Bias and MASE of the estimated functions in the time-invariant (Setting 1) and time-varying (Setting 2) analyses

| Setting 1 | Nonparametric§ | | | | Misspecified† | | | Correctly specified‡ | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Function | Max Relative bias[1] | Mean Bias[2] | Mean Var[3] | MASE($\hat{f}$) | Max Relative Bias | Mean Bias | Mean Var | Max Relative bias | Mean Bias | Mean Var |
| $\hat{\mu}(t)$ | 0.001 | 0.057 | 0.45 | 0.46 | 0.044 | 1.001 | 0.256 | 0.003 | 0.011 | 0.25 |
| $\widetilde{h_g^2}(t)$ | 0.07 | 0.007 | 0.004 | 0.0032 | 0.18 | 0.01 | 0.0038 | 0.134 | 0.008 | 0.003 |
| $\widetilde{h_T^2}(t)$ | 0.04 | 0.013 | 0.005 | 0.0041 | 0.19 | 0.02 | 0.0049 | 0.108 | 0.012 | 0.0043 |

| Setting 2 | Nonparametric | | | | Misspecified | | | Correctly specified | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Function | Max Relative bias | Mean Bias | Mean Var | MASE($\hat{f}$) | Max Relative Bias | Mean Bias | Mean Var | Max Relative bias | Mean Bias | Mean Var |
| $\hat{\mu}(t)$ | 0.001 | 0.056 | 0.48 | 0.49 | 1.01 | 0.91 | 0.31 | 0.0002 | 0.012 | 0.25 |
| $\hat{\beta}_g(t)$ | 0.007 | 0.036 | 0.638 | 0.92 | 0.45 | 0.23 | 0.36 | 0.007 | 0.0024 | 0.18 |
| $\widetilde{h_g^2}(t)$ | 0.02 | 0.004 | 0.005 | 0.0043 | 0.42 | 0.34 | 0.30 | 0.012 | 0.0038 | 0.004 |
| $\widetilde{h_T^2}(t)$ | 0.01 | 0.006 | 0.005 | 0.0049 | 0.51 | 0.37 | 0.36 | 0.001 | 0.0023 | 0.003 |

[1] Max relative bias is defined as: $\max_{t_k \in T_K} |\text{Mean}(\hat{t}(t_k)) - f(t_k)|/f(t_k)$, where $T_K$ is a set of grid points and the average is taken over all repetitions of simulation.

[2] Mean bias is defined as: $\frac{1}{K}\sum_{t_k \in T_k} \text{Mean}(\widehat{f}(t_k)) - f(t_k))$, where $K$ is cardinality of $T_K$ and the average is taken over all repetitions of simulation.

[3] Mean empirical variance is defined as: $\frac{1}{K}\sum_{t_k \in T_k} Var(\widehat{f}(t_k) - f(t_k))$, where $K$ is cardinality of $T_K$ and the variance is taken over all repetitions of simulation.

§ $\mu(t)$ estimated nonparametrically by penalized splines.

† $\mu(t)$ misspecified as $\mu(t) = a_0 + a_1 t + a_2 t^2$.

‡ $\mu(t)$ estimated in a non-linear mixed effects model with a correctly specified functional form.

**Table 2**

Power for testing $\beta_g = 0$ assuming $\mu(t)$ to be a nonparametric function, misspecified parametric functions and a correctly specified nonlinear function ($\alpha$ level =0.05)

| $\mu(t)$ | Analysis | $\beta_g$ | Nonparametric§ | Misspecified: Linear† | Misspecified: Quadratic‡ | Correctly specified* |
|---|---|---|---|---|---|---|
| Logarithm[a] | Type I err | 0 | 0.048 | 0.012 | 0.024 | 0.048 |
| Logarithm | Power | 0.5 | 0.51 | 0.06 | 0.46 | 0.54 |
| Logarithm | Power | 0.75 | 0.79 | 0.2 | 0.76 | 0.79 |
| Logarithm | Power | 1 | 0.94 | 0.53 | 0.93 | 0.96 |
| Gaussian[b] | Type I err | 0 | 0.046 | 0.85 | 0.004 | 0.058 |
| Gaussian | Power | 0.5 | 0.49 | - | 0 | 0.53 |
| Gaussian | Power | 0.75 | 0.93 | - | 0 | 0.93 |
| Gaussian | Power | 1 | 0.94 | - | 0 | 0.99 |
| Sine[c] | Type I err | 0 | 0.045 | 0.99 | 0.99 | 0.048 |
| Sine | Power | 0.5 | 0.46 | - | - | 0.49 |
| Sine | Power | 0.75 | 0.84 | - | - | 0.86 |
| Sine | Power | 1 | 0.95 | - | - | 0.96 |

§ $\mu(t)$ estimated nonparametrically by penalized splines.

† $\mu(t)$ misspecified as $\mu(t) = a_0 + a_1 t$.

‡ $\mu(t)$ misspecified as $\mu(t) = a_0 + a_1 t + a_2 t^2$.

* $\mu(t)$ estimated in a non-linear mixed effects model with a correctly specified functional form.

[a] True $\mu(t) = -34.2 + 81.7\log(0.3(t + 21.7))$

[b] True $\mu(t) = 200 \exp(-0.002(t - 39)^2))$

[c] True $\mu(t) = 150 + 50 \sin(0.2t)$

**Table 3**

Type I error of the permutation test and the misspecified parametric analyses for testing $\beta_g(t) = 0$ in model (4)

| $a$ level | Nonparametric[§] | Misspecified: Linear[†] | Misspecified: Quadratic[‡] |
|---|---|---|---|
| 0.005 | 0.0054 | 0.0045 | 0.0055 |
| 0.01 | 0.0128 | 0.01 | 0.008 |
| 0.05 | 0.0488 | 0.0515 | 0.0505 |
| 0.1 | 0.0914 | 0.1 | 0.1015 |

[§] $\mu(t)$ estimated nonparametrically by penalized splines.

[†] $\mu(t)$ misspecified as $\mu(t) = a_0 + a_1 t$.

[‡] $\mu(t)$ misspecified as $\mu(t) = a_0 + a_1 t + a_2 t^2$.

**Table 4**

Power for testing $\beta_g(t) = 0$ assuming $\beta_g(t)$ to be a nonparametric function, misspecified parametric functions and a correctly specified nonlinear function[*] ($\alpha$ level =0.05)

| $\beta_g(t)$ | Mean $\beta_g(t)$ over $t$ | Nonparametric[§] | Misspecified: Linear[†] | Misspecified: Quadratic[‡] | Correctly specified[*] |
|---|---|---|---|---|---|
| log(0.2t)/10 − 0.2 | −0.004 | 0.19 | 0.19 | 0.08 | 0.19 |
| log(0.5t)/10 + 0.2 | 0.49 | 0.39 | 0.39 | 0.34 | 0.39 |
| 0.8 + 0.1log(0.5t) | 1.09 | 0.98 | 0.98 | 0.98 | 0.98 |
| 3 exp(−0.075(t − 39)²) − 0.5 | −0.001 | 0.73 | 0.04 | 0.42 | 0.99 |
| 0.9 exp(−0.025(t − 39)²) + 0.2 | 0.45 | 0.35 | 0.34 | 0.37 | 0.77 |
| 1.5 exp(−0.075(t − 39)²) + 0.6 | 0.85 | 0.86 | 0.85 | 0.85 | 0.99 |
| 0.85 sin(0.2t) + 0.02 | 0.01 | 0.60 | 0.06 | 0.05 | 0.94 |
| 0.85 sin(0.2t) + 0.5 | 0.49 | 0.51 | 0.36 | 0.3 | 0.86 |
| 0.85 sin(0.2t) + 0.85 | 0.84 | 0.87 | 0.81 | 0.78 | 0.96 |

[*] In all analyses $\mu(t)$ was estimated nonparametrically by penalized splines.

[†] $\beta_g(t)$ misspecified as $\beta_g(t) = \alpha_0 + \alpha_1 t$.

[‡] $\beta_g(t)$ misspecified as $\beta_g(t) = \alpha_0 + \alpha_1 t + \alpha_2 t^2$.

[†] $\beta_g(t)$ estimated in a non-linear mixed effects model with a correctly specified functional form.

**Table 5**

Top ranking SNPs in the time-invariant and time-varying analyses with FHS data[*]

| Trait | SNP | Analysis | Chr | Location | Gene | MAF | LRT** | $p$ value** | LRT-lin§ | $p$-lin§ |
|---|---|---|---|---|---|---|---|---|---|---|
| SBP | rs11065951 | Invariant† | 12 | 110479861 | ATXN2 | 0.052 | 12.36 | $4.4 \times 10^{-4}$ | - | - |
| DBP | rs1052501 | Varying‡ | 3 | 41900402 | ULK4 | 0.192 | 16.50 | $1.0 \times 10^{-4}$ | 8.87 | 0.01 |
| DBP | rs7648578 | Varying | 3 | 41833735 | ULK4 | 0.187 | 16.81 | $9.8 \times 10^{-5}$ | 9.38 | 0.01 |
| DBP | rs2128834 | Varying | 3 | 41837649 | ULK4 | 0.187 | 16.30 | $1.2 \times 10^{-4}$ | 8.57 | 0.01 |
| DBP | rs3774372 | Varying | 3 | 41852418 | ULK4 | 0.183 | 16.29 | $1.2 \times 10^{-4}$ | 7.73 | 0.02 |
| SBP | rs10858911 | Varying | 12 | 88487272 | - | 0.396 | 24.56 | $1.0 \times 10^{-5}$ | 2.87 | 0.24 |
| SBP | rs4757448 | Varying | 11 | 16954369 | PLEKHA7 | 0.340 | 40.12 | $1.0 \times 10^{-6}$ | 1.23 | 0.54 |
| SBP | rs17700056 | Varying | 11 | 16975383 | PLEKHA7 | 0.341 | 38.39 | $1.0 \times 10^{-6}$ | 1.39 | 0.50 |
| SBP | rs7943587 | Varying | 11 | 16812381 | PLEKHA7 | 0.373 | 32.41 | $5.0 \times 10^{-6}$ | 0.49 | 0.78 |
| SBP | rs7121911 | Varying | 11 | 16977903 | PLEKHA7 | 0.208 | 24.75 | $1.8 \times 10^{-5}$ | 24.75 | $4.2 \times 10^{-6}$ |

[*] Significance level $4.81 \times 10^{-4}$ adjusting for multiple comparisons of 265 SNPs by Gao et al. (2008).

[**] Likelihood ratio statistic and $p$ value in a time-varying analysis treating genetic effect as a nonparametric function and estimated by linear spline.

[†] Treating genetic effect as time invariant.

[‡] Treating genetic effect as time varying.

[§] Likelihood ratio statistic and $p$ value in a time-varying analysis treating genetic effect as a parametric linear function.