

# USING SMOOTHING SPLINE ANOVA TO EXAMINE THE RELATION OF RISK FACTORS TO THE INCIDENCE AND PROGRESSION OF DIABETIC RETINOPATHY

YUEDONG WANG<sup>1\*</sup>, GRACE WAHBA<sup>2</sup>, CHONG GU<sup>3</sup>, RONALD KLEIN<sup>4</sup>  
AND BARBARA KLEIN<sup>4</sup>

<sup>1</sup> *Department of Biostatistics, University of Michigan, 1420 Washington Heights, Ann Arbor, MI 48109, U.S.A.*

<sup>2</sup> *Department of Statistics, University of Wisconsin, 1210 W. Dayton St. Madison, WI 53706, U.S.A.*

<sup>3</sup> *Department of Statistics, University of Michigan, Mason Hall, Ann Arbor, MI 58109, U.S.A.*

<sup>4</sup> *Department of Ophthalmology, University of Wisconsin, 610 Walnut St. Madison, WI 53706, U.S.A.*

## SUMMARY

Smoothing spline ANOVA (ANalysis Of VAriance) methods provide a flexible alternative to the standard parametric GLIM (generalized linear models) methods for analysing the relationship of predictor variables to outcomes with data from large epidemiologic studies. These methods allow the visualization of relationships not readily fit by simple GLIM models, and provide for the ability to visualize interactions between the variables. At the same time, they reduce to GLIM models if the data suggest that the added flexibility is unwarranted. Using this method, we investigate risk factors for incidence and progression of diabetic retinopathy in a group of patients with older onset diabetes from the Wisconsin Epidemiological Study of Diabetic Retinopathy. We carry out four analyses to illustrate various properties of this class of methods. Some of the results confirm previous findings with use of standard methods, while others allow the visualization of more complex relationships not evident from the application of parametric methods. © 1997 by John Wiley & Sons, Ltd.

*Statist. Med.*, **16**, 1357–1376 (1997)

No. of Figures: 10    No. of Tables: 1    No. of References: 38

## 1. INTRODUCTION

Many demographic medical studies collect data of the form  $\{y_i, t(i), i = 1, \dots, n\}$ , where  $i$  indexes the  $i$ th study participant,  $y_i$  is 1 or 0, indicating whether a medical condition of interest is present or absent at follow-up, and  $t(i) = (t_1(i), \dots, t_d(i))$  is a vector of  $d$  predictor variables at baseline, that may or may not relate to the likelihood that  $y_i$  is a 1. From these data we wish to estimate  $p(t) = \text{Prob}\{y = 1 | t\}$ , the probability that a person with predictor vector  $t$  presents with the condition of interest at follow-up. We use such estimates for the prevalence of the condition of interest in general populations and to study the sensitivity of  $p$  to the predictor variables, or, if

\* Correspondence to: Y. Wang, Department of Biostatistics, University of Michigan, 1420 Washington Heights, Ann Arbor, MI 48109, U.S.A.

Contract grant sponsor: NIH

Contract grant number: EY09946, P60 DK20572, P30 HD18258, U54 HD 29184, EY03083

Contract grant sponsor: NSF

Contract grant number: DMS9121003, DMS9301511

CCC 0277–6715/97/121357–20\$17.50

© 1997 by John Wiley & Sons, Ltd.

*Received December 1995*

*Revised August 1996*

possible, to combinations of them. The traditional GLIM models<sup>1</sup> do this by defining  $f(t)$ , the logit, as

$$f(t) = \log[p(t)/(1 - p(t))] \quad (1)$$

and assuming that  $f$  is a simple parametric function of the components of  $t$ . When the  $t_x$  are continuous variables, the most commonly used model is linear in the components of  $t$

$$f(t) = f(t_1, \dots, t_d) = \mu + \sum_{\alpha=1}^d \beta_{\alpha} t_{\alpha} \quad (2)$$

but sometimes one uses second or even third degree polynomials if it appears that a linear model is inadequate. If some of the  $t_x$  are categorical variables, then we can use indicator functions. More generally, given  $M$  parametric functions  $\phi_v$ ,  $v = 1, \dots, M$  subject to some identifiability conditions,  $f$  is modelled parametrically as

$$f(t) = \sum_{v=1}^M \beta_v \phi_v(t). \quad (3)$$

Then ( $\mu$  and) the  $\beta$ 's are obtained by minimizing the negative log-likelihood  $\mathcal{L}(y, f)$  given by

$$\mathcal{L}(y, f) = \sum_{i=1}^n [y_i f(t(i)) - \log(1 + e^{f(t(i))})] \quad (4)$$

with (2) or, more generally, (3), substituted for  $f$ . If the 'true but unknown'  $f$  is actually some linear combination of the specified  $\phi_v$  then this is, of course the correct procedure. Unfortunately, as we examine large data sets more closely, it becomes clear that linear models, or even quadratic or cubic models, are often inadequate. What happens if, for example, the dependence on  $t_x$  is 'J' shaped, or, even has two peaks, possibly representing two distinct subpopulations? If the 'true' log odds ratio  $f$  is not well approximated by some function of the specified form, then we may have large biases in the estimates of  $f$ . Furthermore, the common statistical hypothesis tests, confidence intervals,  $P$ -values, etc. are not necessarily valid if  $f$  is not of the specified form. In this paper we describe and demonstrate the use of smoothing spline ANOVA (SS-ANOVA) methods to estimate  $f$ . These methods avoid many of the above difficulties, yet they provide the ability to visualize some of the relationships between the variables not easily observed with use of more traditional methods. We demonstrate their use to examine the relation of risk factors to the incidence and progression of diabetic retinopathy, using data from the Wisconsin Epidemiological Study of Diabetic Retinopathy (WESDR).

The analysis here is based on (a special case of) the smoothing spline ANOVA method for modelling and estimating  $f$  that appears in Wahba, Wang, Gu, Klein and Klein<sup>2</sup> (WWGKK), and we have implemented the analysis via the publicly available code GRKPACK described in Wang.<sup>3</sup> The code itself is available from [netlib@research.att.com](mailto:netlib@research.att.com) in the `gcv` directory there, and through [statlib@lib.stat.cmu.edu](mailto:statlib@lib.stat.cmu.edu).

An SS-ANOVA estimate is a form of penalized likelihood estimate. We first note two important penalized likelihood estimates that appear in the literature. Then, in the remainder of this introduction, we describe the main features of SS-ANOVA in general and the details of the SS-ANOVA models in particular that we apply to the WESDR data.

Major precursors of the work under discussion include O'Sullivan<sup>4</sup> and O'Sullivan *et al.*,<sup>5</sup> who proposed a penalized log-likelihood estimate for  $f$  based on thin plate splines. These splines are useful in many contexts, and one can incorporate them in an SS-ANOVA model;<sup>6,7</sup> we do not employ them in the present work.

Hastie and Tibshirani<sup>8</sup> (see also other references there) discussed estimates of  $f$  of the form

$$f(t) = f(t_1, \dots, t_d) = \mu + \sum_{\alpha=1}^d f_{\alpha}(t_{\alpha}) \quad (5)$$

where the  $f_{\alpha}$  are 'smooth' functions obtained by some smoothing process, including the use of cubic smoothing splines. The S code (Chambers and Hastie<sup>9</sup>) provides the facility for fitting models of the form (5). They also note that some of their work extends to the SS-ANOVA methods that we consider here. The varying-coefficient models of Hastie and Tibshirani<sup>10</sup> are a very interesting sub-family of the SS-ANOVA models. Cubic smoothing splines for the  $f_{\alpha}$  in (5) are the solution to the minimization problem: find  $f_{\alpha}$  (in an appropriate function space), and subject to some identifiability criteria such as  $\int f_{\alpha}(t_{\alpha}) dt_{\alpha} = 0$ , to minimize the penalized log-likelihood functional

$$\mathcal{J}_{\lambda}(y, f) = \mathcal{L}(y, f) + \sum_{\alpha=1}^d \lambda_{\alpha} J_{\alpha}(f_{\alpha}) \quad (6)$$

where we define the penalty functionals  $J_{\alpha}$  by

$$J_{\alpha}(f_{\alpha}) = \int_0^1 (f_{\alpha}''(t_{\alpha}))^2 dt_{\alpha}. \quad (7)$$

As the smoothing parameter  $\lambda_{\alpha}$  tends to infinity,  $f_{\alpha}$  tends to a linear function in  $t_{\alpha}$ , so that the minimizer of (6) tends to the linear function as in (2) if all the  $\lambda_{\alpha}$ 's become large. Penalized likelihood methods with more general penalty functionals and unpenalized terms are discussed in reference 11.

Recent research has focused on more general models for  $f$  that allow the explicit modelling and visualization of possible interactions between variables, via functional analysis of variance decompositions (to be described) and SS-ANOVA estimation methods. Given a fairly arbitrary function  $f(t_1, \dots, t_d)$  of several variables, we can define a (functional) ANOVA decomposition of  $f$  (generalizing ideas from parametric ANOVA) as

$$f(t_1, \dots, t_d) = \mu + \sum_{\alpha=1}^d f_{\alpha}(t_{\alpha}) + \sum_{\alpha < \beta} f_{\alpha\beta}(t_{\alpha}, t_{\beta}) + \dots + f_{1, \dots, d}(t_1, \dots, t_d) \quad (8)$$

where the  $f_{\alpha}$  are the main effects,  $f_{\alpha\beta}$  are the two factor interactions, etc. The components are uniquely determined given a set of averaging operators  $\mathcal{E}_{\alpha}$ , which average functions over the  $t_{\alpha}$  in some specified way. For example, if  $t_{\alpha}$  is a continuous variable in the interval  $[0, 1]$ , then a possible choice for  $\mathcal{E}_{\alpha}$  is

$$(\mathcal{E}_{\alpha} f)(t_1, \dots, t_{\alpha-1}, t_{\alpha+1}, \dots, t_d) = \int_0^1 f(t_1, \dots, t_d) dt_{\alpha}. \quad (9)$$

Then the mean is  $\mu = \prod_{\alpha=1}^d \mathcal{E}_{\alpha} f$ , the main effects are  $f_{\alpha} = (I - \mathcal{E}_{\alpha}) \prod_{\beta \neq \alpha} \mathcal{E}_{\beta} f$ , the two factor interactions are  $f_{\alpha\beta} = (I - \mathcal{E}_{\alpha})(I - \mathcal{E}_{\beta}) \prod_{\gamma = \alpha, \beta} \mathcal{E}_{\gamma} f$ , etc. Since  $\mathcal{E}_{\alpha} \mathcal{E}_{\alpha} = \mathcal{E}_{\alpha}$ , we can see that the terms in the functional ANOVA decomposition satisfy side conditions analogous to those in ordinary parametric ANOVA, for example,  $\mathcal{E}_{\alpha} f_{\alpha} = 0$ . In general, for model fitting purposes, we eliminate higher order terms in the functional ANOVA decomposition and we estimate some or all of the lower order terms by finding  $\mu$ , and (some of) the  $f_{\alpha}, f_{\alpha\beta}$ , etc. to minimize a penalized log-likelihood functional  $\mathcal{J}_{\lambda}(y, f)$  given by an expression that generalizes (6), with a separate penalty functional for each independently smoothed term in the model. We can also incorporate in the

model indicator functions for categorical variables. Other averaging operators include weighted sums over possible or observed values of  $t_x$ .

It is well known<sup>11</sup> that solutions to variational problems like the minimizer of  $\mathcal{J}_\lambda(y, f)$  and generalizations to  $\mathcal{J}_\lambda(y, f)$  with  $f$  containing interaction terms as in (8) have a representation as a linear combination of the unpenalized terms in the model plus  $n$  (basis) functions, that we can construct from the  $t(i)$ , the reproducing kernels associated with the penalty functionals, and the smoothing parameters, see references 3, 4, 11–14. Given the smoothing parameters, the numerical problem of computing the minimizer  $f_\lambda$  reduces to finding the coefficients of a representation for it as a linear combination of the unpenalized terms and the above mentioned basis functions. We can solve the numerical problem for the coefficients with a Newton–Raphson iteration in conjunction with matrix decompositions. The history of these numerical methods includes references 4, 5, 11, 15–18. There are various approximate numerical methods available and under development for data sets that are too large for matrix decomposition methods; these, however, are beyond the scope of this article.

Objective methods for choosing the smoothing parameters are desirable. This paper and the code GRKPACK employ the iterative unbiased risk method given in Gu<sup>15,16</sup> for the non-Gaussian case and extended to the  $p$  smoothing parameter case  $\lambda = (\lambda_1, \dots, \lambda_p)$  in Wang<sup>19</sup> and WWGKK. This method is a computable proxy for the Kullback–Liebler information distance from the estimate to the ‘truth’, see references 2, 16, 20 and 21. Thus, the method attempts to choose the smoothing parameters to minimize the Kullback–Liebler information distance from the estimate  $f_\lambda$  and the unknown ‘true’  $f$ . Other related references are 22–24.

As with any estimation method, it is important to have some indication of its accuracy. It is particularly important in the case here of non-parametric regression with Bernoulli data derived from medical records because these data tend to be irregular, and may contain influential outliers. The rigidity of parametric models may mask the effect of outliers, as well as obscure the fact that we sometimes make inferences in fairly data-sparse regions. The non-parametric estimates may be more sensitive to outliers, and it is important to be able to delineate a region in the predictor variable space where we can rely upon the estimate, as well as where we can provide some sort of confidence statement. In this work we use the Bayesian ‘confidence intervals’ proposed in reference 25, adapted to the component-wise multiple smoothing parameter case in reference 7, to non-Gaussian data in reference 26, and to the multiple smoothing parameter non-Gaussian case in reference 19 and WWGKK. We use the Bayesian ‘confidence interval’ for  $f_\lambda$  to delineate the region in predictor variable space for which we deem the overall estimate as reliable, by computing an appropriate level curve in a contour plot of the width of the confidence interval, and using this to enclose a ‘reliable’ region. We also use it in conjunction with cross-sectional plots so that we can visualize the accuracy estimates. We have used the component-wise confidence intervals to eliminate some terms that we cannot distinguish from noise, by deleting terms whose confidence intervals contain 0 over most of their domain. We remark that we base these confidence intervals on an across-the-function property, that is, a 95 per cent confidence interval has the property that (approximately) we expect the confidence intervals to cover the true curve at about 95 per cent of the  $n$  data points, see references 3, 21 and 26.

We now proceed to the details of the particular models employed in the analysis of the WESDR data. We carry out four analyses, each of which illustrates some particular feature of this class of models. In each case we have rescaled the (continuous) predictor variables  $t_1, \dots, t_d$  to  $[0, 1]$ , and we use the averaging operator  $\mathcal{E} = \mathcal{E}_x$  defined in (9). We consider only main effects and two factor interactions. We employ the penalty functional  $J_x$  of (7) and a two-dimensional relative  $J_{x\beta}$  defined below. If we eliminated all of the two factor interactions, then the models described

immediately below reduce to the form (5) as studied by Hastie and Tibshirani,<sup>8</sup> although our estimation method is different.

We require some definitions to define the terms in the particular ANOVA decomposition (8) that we use. Define the (linear) ‘trend’ function  $\phi(u) = u - 1/2$  for  $u \in [0, 1]$ , and, for continuous functions  $g$  defined on  $[0, 1]$  let  $\mathcal{B}g$  stand for  $\mathcal{B}g = g(1) - g(0)$ . For future reference note that  $\mathcal{E}\phi \equiv \int_0^1 \phi(u) du = 0$  and  $\mathcal{B}\phi = 1$ . We use  $\mathcal{E}$  and  $\mathcal{B}$  to define ANOVA and identifiability side conditions. A subscript  $\alpha$  on  $\mathcal{E}$  or  $\mathcal{B}$  means that it applies to what follows considered as a function of  $t_\alpha$ . The typical main effect term  $f_\alpha(t_\alpha)$  in this set up has a decomposition of the form

$$f_\alpha(t_\alpha) = d_\alpha \phi(t_\alpha) + f_{s\alpha}(t_\alpha) \quad (10)$$

where  $d_\alpha \phi(t_\alpha)$  is the linear and unpenalized part of  $f_\alpha$  and  $f_{s\alpha}$  is the (detrended) ‘smooth’ part of  $f_\alpha$ .  $f_{s\alpha}$  satisfies the side conditions  $\mathcal{E}_\alpha f_{s\alpha} = \mathcal{B}_\alpha f_{s\alpha} = 0$  and appears inside a penalty functional as  $J_\alpha(f_{s\alpha})$ . Due to the side conditions on  $f_\alpha$ ,  $J_\alpha(f_{s\alpha}) = 0$  implies that  $f_{s\alpha} = 0$ .

The two factor interaction  $f_{\alpha\beta}(t_\alpha, t_\beta)$  has an analogous decomposition into four interaction terms, namely

$$f_{\alpha\beta}(t_\alpha, t_\beta) = d_{\alpha\beta} \phi(t_\alpha) \phi(t_\beta) + \phi(t_\alpha) f_{s\beta}^{(\alpha)}(t_\beta) + \phi(t_\beta) f_{s\alpha}^{(\beta)}(t_\alpha) + f_{s\alpha\beta}(t_\alpha, t_\beta) \quad (11)$$

where the ‘smooth’ factors of the trend  $\times$  smooth interactions satisfy the side conditions  $\mathcal{E}_\alpha f_{s\alpha}^{(\beta)} = \mathcal{B}_\alpha f_{s\alpha}^{(\beta)} = \mathcal{E}_\beta f_{s\beta}^{(\alpha)} = \mathcal{B}_\beta f_{s\beta}^{(\alpha)} = 0$  and the ‘smooth–smooth’ interaction term satisfies  $(\mathcal{E}_\alpha f_{s\alpha\beta})(t_\beta) = (\mathcal{B}_\alpha f_{s\alpha\beta})(t_\beta) = (\mathcal{E}_\beta f_{s\alpha\beta})(t_\alpha) = (\mathcal{B}_\beta f_{s\alpha\beta})(t_\alpha) = 0$ , for all  $t_\alpha, t_\beta$ . Letting  $J_{\alpha\beta}(f_{s\alpha\beta}) = \int_0^1 \int_0^1 \left( \frac{\partial^2}{\partial t_\alpha^2 \partial t_\beta^2} f_{s\alpha\beta}(t_\alpha, t_\beta) \right)^2 dt_\alpha dt_\beta$ , then one can show that the side conditions insure that  $J_{\alpha\beta}(f_{s\alpha\beta}) = 0$  implies that  $f_{s\alpha\beta} = 0$ . Letting  $\lambda = (\lambda_\alpha, \lambda_\beta, \lambda_\alpha^{(\beta)}, \lambda_\beta^{(\alpha)}, \lambda_{\alpha\beta})$  and

$$\mathcal{J}_\lambda(f) = \lambda_\alpha J_\alpha(f_{s\alpha}) + \lambda_\beta J_\beta(f_{s\beta}) + \lambda_\alpha^{(\beta)} J_\alpha(f_{s\alpha}^{(\beta)}) + \lambda_\beta^{(\alpha)} J_\beta(f_{s\beta}^{(\alpha)}) + \lambda_{\alpha\beta} J_{\alpha\beta}(f_{s\alpha\beta}) \quad (12)$$

we estimate  $f$  by finding  $\mu, d_\alpha, d_\beta, d_{\alpha\beta}$ , and  $f_{s\alpha}, f_{s\beta}, f_{s\alpha}^{(\beta)}, f_{s\beta}^{(\alpha)}$  and  $f_{s\alpha\beta}$  to minimize

$$\mathcal{J}_\lambda(y, f) = \mathcal{L}_\lambda(y, f) + \mathcal{J}_\lambda(f). \quad (13)$$

Due to the side conditions, as a component of  $\lambda$  becomes large, the estimate of the ‘smooth’ term that it multiplies becomes small. Thus, a good method for choosing the components of  $\lambda$  from the data effectively eliminates unneeded terms in the expansion if one estimates their companion smoothing parameters as large.

We can include categorical variables in the model, by, for example, letting

$$f(t_\alpha, t_\beta, z) = \mu + f_\alpha(t_\alpha) + f_\beta(t_\beta) + f_{\alpha\beta}(t_\alpha, t_\beta) + \sum_{k=2}^K \gamma_k I_k(z) \quad (14)$$

where  $z$  is a variable with  $K$  possible values  $z_1, \dots, z_K$  and  $I_k(z) = 1$  if  $z = z_k$  and 0 otherwise. We then include the  $\gamma$ 's in the minimization of (13).

The mathematics behind the representation of  $f$ , the numerical method for the minimization of  $\mathcal{J}_\lambda$ , the data-based choice of  $\lambda$  according to the iterative unbiased risk method, and the calculation of the confidence intervals are described in WWGKK. Further details appear in reference 3 and the documentation for GRKPACK. In the four analyses that we carry out below, there are between  $d = 2$  and  $d = 4$  predictor variables. In the  $d = 4$  case, considering the penalty functional in (12), and including all of the possible main effects and two factor interaction terms with their own smoothing parameters, the result is 4 main effects smoothing parameters ( $\lambda_\alpha$ ), 12 trend  $\times$  smooth smoothing parameters ( $\lambda_\alpha^{(\beta)}$ ) and 6 smooth  $\times$  smooth parameters ( $\lambda_{\alpha\beta}$ ), more than we would want to fit simultaneously with the sample sizes here (all less than 500). We reduced the

number of terms with smoothing parameters to a maximum of 7 based on previously published analyses of these data by more traditional methods, by extensive pre-screening, by fitting parametric polynomial GLIM models with terms as high as cubic order to detect and delete terms that failed to give any evidence of significance, and by fitting some marginal SS-ANOVA models as demonstrated in Section 6. Exploration of more systematic pre-screening methods is an area of active research. Once we had reduced the number of smoothing parameters to 7 or less, we used informal model selection methods as discussed in WWGKK. These included the deletion of component terms too small to have an observable effect on cross-sectional plots of  $f$ , and the deletion of terms whose componentwise confidence intervals included the 0 function. Further details specific to each analysis appear below. Leaving-out-one-third model selection procedures have been discussed in reference 20, but we have not used them here.

Section 2 discusses the Wisconsin Epidemiologic Study of Diabetic Retinopathy, and Sections 3, 4, 5 and 6 discuss four analyses of this study. Section 7 is a summary and conclusion.

## 2. WISCONSIN EPIDEMIOLOGIC STUDY OF DIABETIC RETINOPATHY (WESDR)

The WESDR is an ongoing epidemiologic study of a cohort of patients who receive their medical care in an 11-county area in southern Wisconsin, who were first examined in 1980–1982, then again in 1984–1986 and 1990–1992. At this writing a third follow-up is currently in progress. Detailed descriptions of the data have appeared in references 27 and 28 and references there. In brief, a sample of 2990 diabetic patients was selected in an 11-county area in southern Wisconsin. This sample consisted of two groups. The first group was 1210 patients diagnosed as diabetic before they were 30 years of age and who took insulin ('younger onset group'). The second group consisted of 1780 patients diagnosed with diabetes after 30 years of age. Of these, 824 were taking insulin ('older onset group taking insulin') and 956 were not ('older onset group not taking insulin'). Of the 2990 eligible patients, 2366 participated in the baseline examination from 1980 to 1982.

A large number of medical, demographic, ocular and other covariates were recorded at the baseline and later examinations along with a retinopathy score for each eye (to be described). Relations between various covariates and the retinopathy scores have had extensive analysis with standard statistical methods including categorical data analysis and linear logistic models, with results reported in a series of WESDR manuscripts 29–35. We applied SS-ANOVA methods to a subset of the data from the younger onset group in WWGKK in conjunction with an extensive technical account of the mathematical theory and numerical methods behind the method. It is the purpose of this account to carry out four SS-ANOVA analyses on data from the older onset WESDR group, and in the process, to illustrate some of the more important features of the method. Our goal is to explain the possible results, advantages and disadvantages in a less technical manner, aimed at clinicians and medical researchers.

We limited the present study to the construction of predictive models for incidence and for progression (to be defined) of diabetic retinopathy at the first follow-up, as a function of some of the covariates available at baseline. We do this both for the 'older onset group taking insulin' (ID group) and the 'older onset group not taking insulin' (NID group). We analyse both 'incident events' and 'progression events' for both the ID and the NID groups. We only list the covariates pertinent to our analysis:

1. Age: age at the examination (years);
2. Duration: duration of diabetes at the examination (years);
3. Glycosylated haemoglobin: a measure of hyperglycaemia (%);<sup>8</sup>

4. Body mass index (bmi): weight in kg/(height in m)<sup>2</sup>;
5. Pulse rate counted for 30 seconds;
6. Baseline retinopathy severity levels (base-retinopathy-level). See explanation below.

At the baseline and follow-up examinations, each eye was graded as one of the 6 levels: 10, 21, 31, 41, 51 and 60+, in order of increasing retinopathy severity with 10 indicating no retinopathy and 60+ indicating the most severe stage, proliferative retinopathy. We derived the retinopathy level for a participant by giving the eye with the higher level (more severe retinopathy) greater weight. For example, we specify the level for a participant with level 31 retinopathy in each eye with the notation 'level 31/31'; for a participant with level 31 in one eye and less severe retinopathy in the other eye the notation is 'level 31/<31'. This scheme provided an 11-step scale: 10/10; 21/<21; 21/21; 31/<31; 31/31; 41/<41; 41/41; 51/<51; 51/51; 60+/<60+, and 60+/60+. Participants in the analysis for *incidence* consisted of subjects with level 10/10 at the baseline and no missing data. The model provides an estimate of the incidence rate, defined as the probability that such a participant has level 21/<21 or worse at the follow-up examination. Participants in the analysis for *progression* consisted of subjects with no or non-proliferative retinopathy at baseline and no missing data. The model provides an estimate of the progression rate, with progression defined as an increase in baseline level by two steps or more (10/10 to 21/21 or greater, or 21/<21 to 31/<31 or greater, for instance). These analyses correspond to analyses previously carried out by traditional methods in some of the WESDR references cited. Note that this allows subjects included in the incidence analysis also to be part of the progression analysis.

### 3. INCIDENCE IN THE OLDER ONSET GROUP NOT TAKING INSULIN

After excluding participants with missing values, there were 297 participants in the older onset NID group with a baseline score of 10/10, thus qualifying them for inclusion in the NID 'Incidence' analysis. There was one observation of glycosylated haemoglobin recorded as 23.6 per cent, much greater than the others. We decided to delete this influential observation. Our conclusions would remain the same if we included this observation in the analysis.

Klein *et al.*<sup>27</sup> found, using linear logistic models, that glycosylated haemoglobin is the only statistically significant predictor of incidence of retinopathy in older onset patients. Using linear logistic models as a screening tool, we found that the effect of age is not linear and that there is a strong interaction between age and glycosylated haemoglobin.

We first fit an SS-ANOVA model with the main effects of age and glycosylated haemoglobin and all three interaction terms. The main effect of glycosylated haemoglobin is linear. All interaction terms except  $trend(\text{glycosylated haemoglobin}) \times smooth(\text{age})$  are near zero. The final model is

$$\begin{aligned}
 f(\text{age, glycosylated haemoglobin}) \\
 = \mu + f_1(\text{age}) + a_1 \times \text{glycosylated haemoglobin} \\
 + trend(\text{glycosylated haemoglobin}) \times smooth(\text{age}). \quad (15)
 \end{aligned}$$

We plot age versus glycosylated haemoglobin on Figure 1(a). We have marked those participants with incident retinopathy as solid circles and those without as open circles. We superimpose the contour lines of estimated posterior standard deviations. These contours agree well with the distribution of the observations. Thus we can use them to delineate a region in which

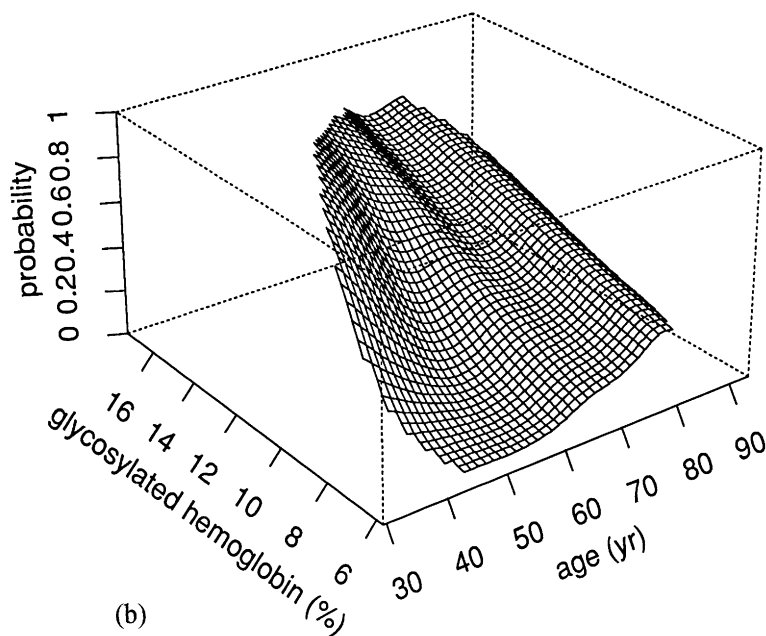
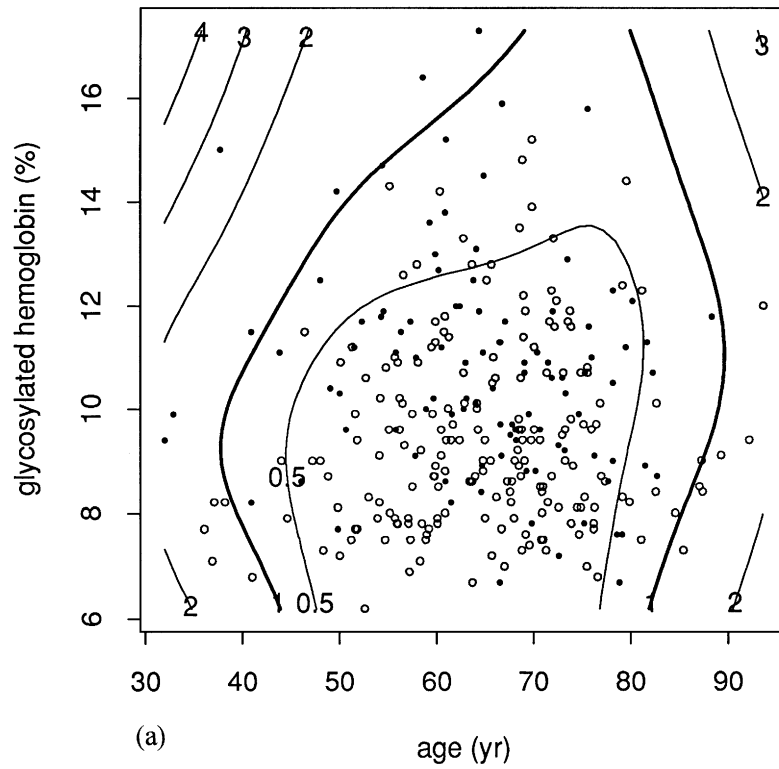


Figure 1. NID 'Incidence' analysis: (a) data and contours of constant posterior standard deviation. Solid circles indicate incidence and open circles indicate non-incidence; (b) estimated probability of incidence in the defined region, as a function of age and glycosylated haemoglobin



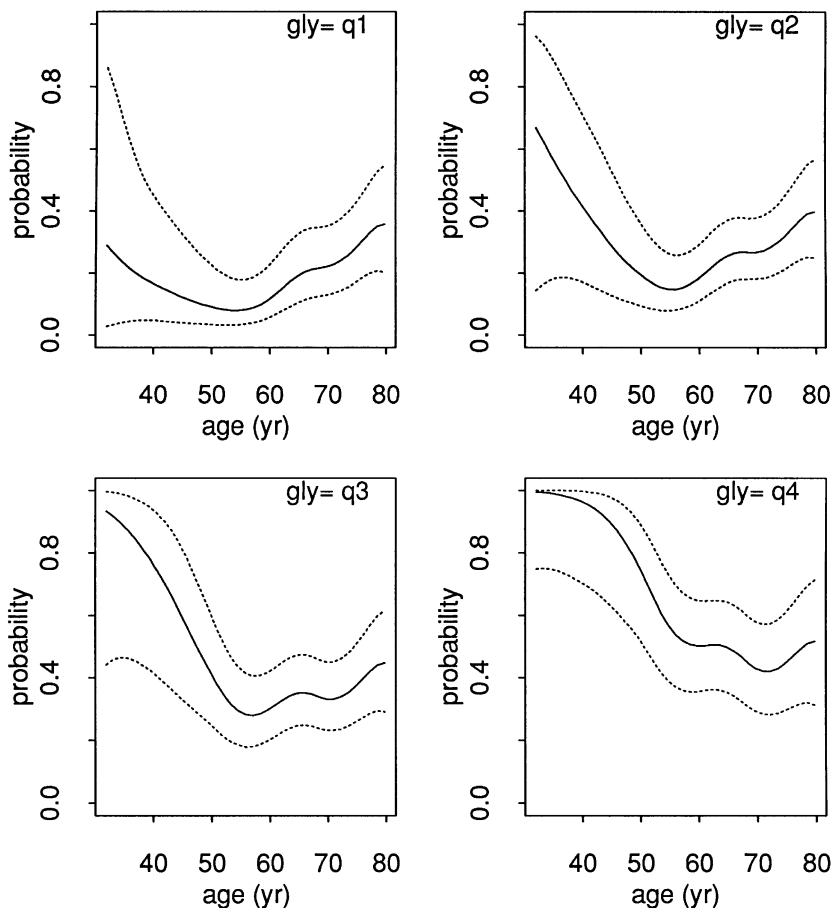


Figure 2. NID 'Incidence' analysis. Cross-sections of estimated probability of incidence as a function of age, with their 90 per cent Bayesian confidence intervals, at four quantiles of glycosylated haemoglobin. q1, q2, q3 and q4 are the quantiles at 0-125, 0-375, 0-625 and 0-875

we deem the estimate of the probability function as reliable. We decided to use the region with estimated posterior standard deviations less than or equal to 1, in the logit scale.

We plot the probability function estimate on Figure 1(b). Figure 2 gives cross-sections of the estimated probability of incidence as a function of age with their 90 per cent Bayesian confidence intervals at the cross-sections, at four quantiles of glycosylated haemoglobin. The width of the confidence intervals suggests that the small bumps are probably an artifact. The general shape of the response, however, is evident, and it appears that the age effect is not particularly well modelled with a second or even a third degree polynomial.

The risk for incidence of retinopathy increases with increasing glycosylated haemoglobin. The risk increases with increasing age for lower glycosylated haemoglobin. This increase might result from the development of other conditions such as hypertension and atherosclerotic vascular disease in older compared to younger subjects and that leads to increased risk of retinopathy, even when the glycosylated haemoglobin level is relatively low. The risk decreases with increasing age for higher glycosylated haemoglobin. This decrease may result from mortality since an old participant with high glycosylated haemoglobin has a higher mortality risk during the observation period.<sup>36</sup>

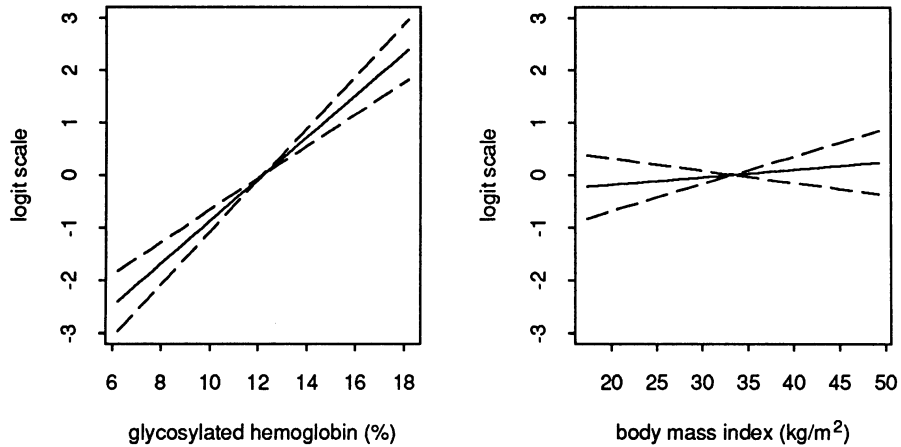


Figure 3. NID 'Progression' analysis. Main effects estimates of glycosylated haemoglobin and bmi along with their component-wise 90 per cent Bayesian confidence intervals, logit scale

#### 4. PROGRESSION OF THE OLDER ONSET GROUP NOT TAKING INSULIN

After excluding participants with missing values, there were 432 participants in the older onset NID group qualified for inclusion in the 'Progression' analysis. Klein *et al.*<sup>27</sup> found that age, duration, and glycosylated haemoglobin were statistically significant in a linear logistic model. Using linear logistic models, we found that the duration main effect of polynomials up to the cubic is significant and the bmi main effect of polynomials up to the quartic is significant. We also found that the multiplication interaction between age and polynomials up to cubic of duration is significant. These indicate that a linear logistic model with lower order polynomials may be inadequate. We decided to fit the SS-ANOVA model:

$$\begin{aligned}
 &f(\text{age, duration, glycosylated haemoglobin, bmi}) \\
 &= \mu + f_1(\text{age}) + f_2(\text{duration}) + f_{1,2}(\text{age, duration}) \\
 &\quad + f_3(\text{glycosylated haemoglobin}) + f_4(\text{bmi}). \tag{16}
 \end{aligned}$$

Figure 3 plots the main effects estimates and their 90 per cent Bayesian confidence intervals for glycosylated haemoglobin and bmi. We see that the effect of glycosylated haemoglobin is strong. The effect of bmi is small, and 0 is contained within its confidence interval. The estimates are effectively linear on a logit scale even though we allowed a 'smooth' term for each of them in the model. This illustrates the ability of the method to reduce to a partially parametric form if warranted by the data.

The effects of glycosylated haemoglobin and bmi are additive in the logit scale. In the remaining plots, we fix glycosylated haemoglobin and bmi at their median values.

We plot age versus duration on Figure 4(a). Participants who had progression of retinopathy appear as solid circles and those with no progression as open circles. We superimpose contour lines of the estimated posterior standard deviations as a function of age and duration with glycosylated haemoglobin and bmi fixed at their median values. These contours agree well with the distribution of the observations. We decided to use the region with estimated

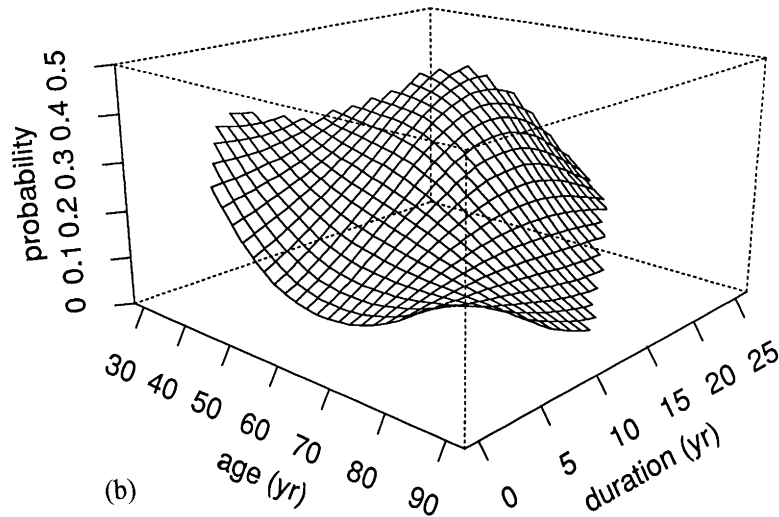
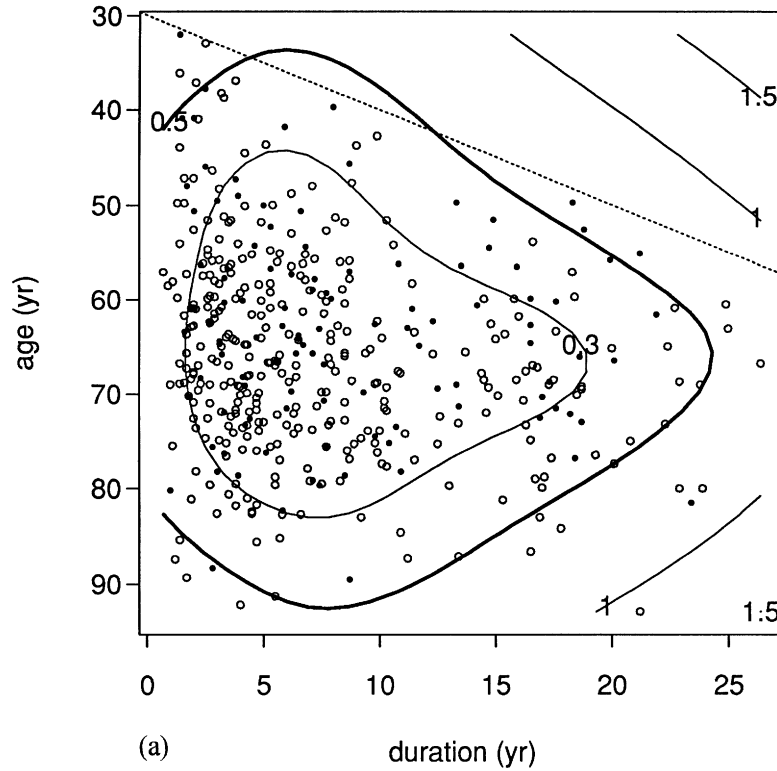


Figure 4. NID 'Progression' analysis: (a) data and contours of constant posterior standard deviation as a function of age and duration at the median value of glycosylated haemoglobin and bmi. The dotted line is age - duration = 30 years; (b) estimated probability of progression in the defined region, as a function of age and duration, at the median value of glycosylated haemoglobin and bmi

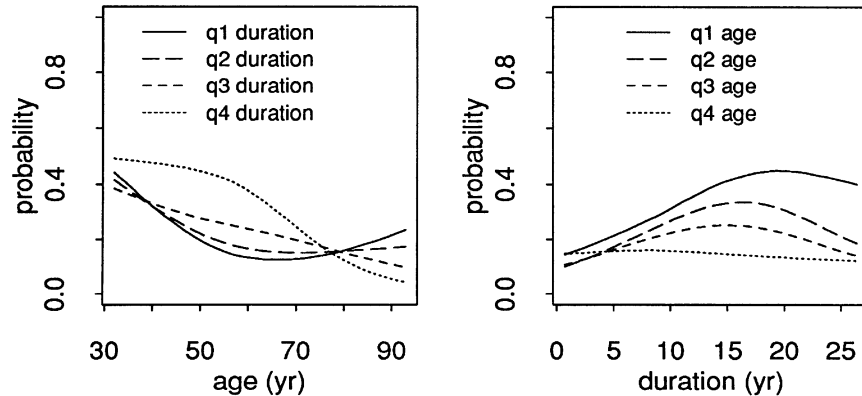


Figure 5. NID 'Progression' analysis. Cross-sections of estimated probability of progression as a function of age and duration, at the median value of glycosylated haemoglobin and bmi

posterior standard deviations less than or equal to 0.5. A plot of the probability function estimate appears on Figure 4(b).

To see more clearly how the probability of progression depends on age and duration, we plot the cross-sections of the estimate in Figure 5. The cross-sections with their 90 percent Bayesian confidence intervals appear in Figures 6 and 7. From these plots, we see that the risk of progression of retinopathy decreases with increasing age and the shapes differ between the fourth and other quantiles of duration. The risk increases with increasing duration up to about 17 years and does not increase after that. Increased mortality of those with longer duration and more severe retinopathy may explain, in part, this finding. Table I gives the correspondence between the percentiles and the physical units.

Plots in Figure 4 indicate that participants diagnosed with diabetes just above 30 years of age (points just below the dotted line in the left panel) are at higher risk of progression of their retinopathy than those diagnosed later in life.

## 5. INCIDENCE IN THE OLDER ONSET GROUP TAKING INSULIN

After excluding participants with missing data, there were 143 participants in the older onset ID group that had a baseline score of 10/10, and hence are included in ID 'Incidence' analysis. This sample size is probably not large enough to estimate interactions well. Klein *et al.*<sup>28</sup> found that age and glycosylated haemoglobin are significant using a linear logistic model. We found that quadratic terms of duration and pulse are significant using a linear logistic:

$\text{logit}(p(\text{age, duration, glycosylated haemoglobin, pulse}))$

$$= \mu + a_1 \times \text{age} + a_2 \times \text{duration} + a_3 \times \text{duration}^2 + a_4 \times \text{glycosylated haemoglobin} + a_5 \times \text{pulse} + a_6 \times \text{pulse}^2. \quad (17)$$

Fitting an SS-ANOVA model with main effects of age, duration, glycosylated haemoglobin and pulse, we found that the main effects of age and glycosylated haemoglobin are linear. Then we fitted the model:

$f(\text{age, duration, glycosylated haemoglobin, pulse})$

$$= \mu + a_1 \times \text{age} + f_1(\text{duration}) + a_2 \times \text{glycosylated haemoglobin} + f_2(\text{pulse}). \quad (18)$$

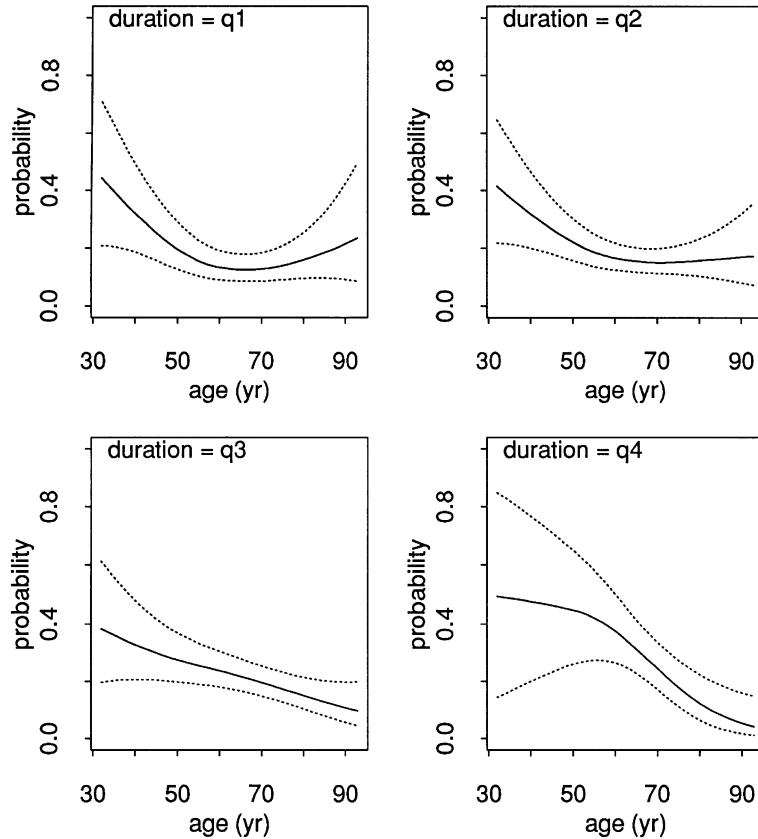


Figure 6. NID 'Progression' analysis. Cross-sections of estimated probability of progression as a function of age with their 90 per cent Bayesian confidence intervals, at four quantiles of duration, at the median values of glycosylated haemoglobin and bmi

The estimates of the main effects and their 90 percent Bayesian confidence intervals appear in Figure 8. We believe the ups and downs in the middle of the main effect of duration are caused by the poor choice of the smoothing parameter (too small) due to the small sample size, but the pattern of the main effect of duration is reliable and agrees with the previous conclusion. That is, the risk increases with increasing duration up to about 6 years and thereafter does not increase any more. From the main effect of pulse, we conclude that the risk is higher for participants with higher pulse rate. Higher resting pulse rate may be secondary to diabetic neuropathy involving the autonomic nervous system, a condition postulated as involved in the development of diabetic retinopathy. The fits from model (17) are well inside the Bayesian confidence intervals of the SS-ANOVA estimates. It is difficult to distinguish between these two models with such a small sample size. The quadratic form of the linear logistic model for pulse is necessarily symmetric about its minimum at around 35 and suggests that the risk is increasing as pulse decreases below 35 while the SS-ANOVA model suggests that the minimum is a little higher, and the curve flattens out below its minimum.

The SS-ANOVA method needs relatively large sample sizes to obtain good estimates of multiple smoothing parameters. Our experience with real data and simulations suggests that about 100 observations for each smoothing parameter gives fairly reliable estimates.

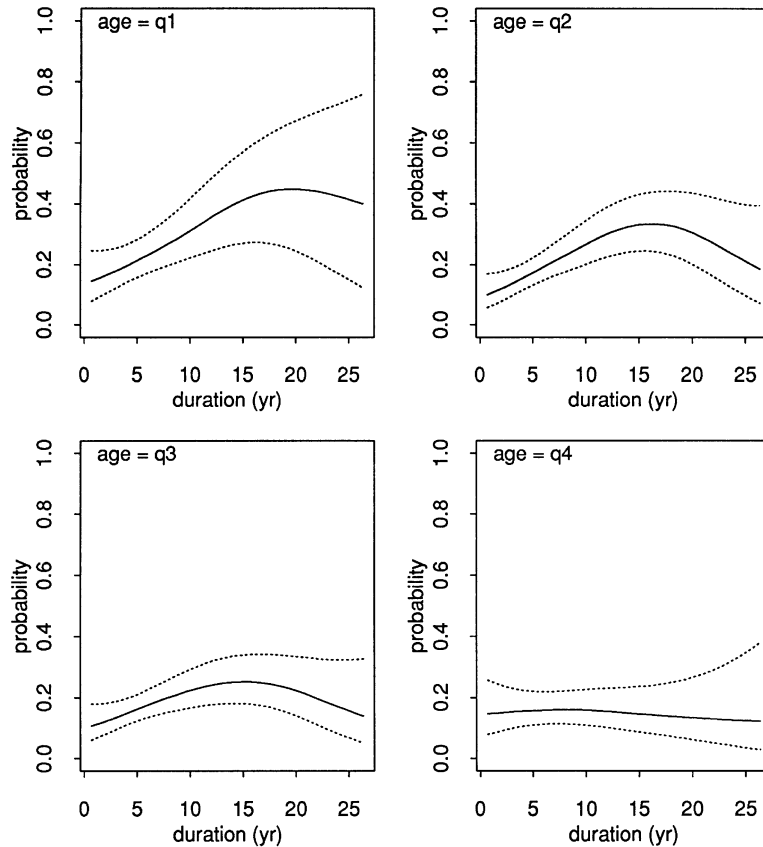


Figure 7. NID 'Progression' analysis. Cross-sections of estimated probability of progression as a function of duration, with their 90 per cent Bayesian confidence intervals, at four quantiles of age and at the median values of glycosylated haemoglobin and bmi

Table I. Percentiles used in plots

| Percentile    | 12.5 | 37.5 | 62.5 | 87.5 |
|---------------|------|------|------|------|
| Age (yr)      | 53.0 | 63.5 | 69.8 | 78.2 |
| Duration (yr) | 2.5  | 4.2  | 7.8  | 16.7 |

## 6. PROGRESSION OF THE OLDER ONSET GROUP TAKING INSULIN

Due to its flexibility, we can use the SS-ANOVA method in several stages of data analyses. In the previous sections, we explained the SS-ANOVA method as a tool for model building and estimation. In practice, we can also use the SS-ANOVA method to explore the behaviour of raw data. For example, we can obtain a marginal estimate for each covariate to investigate whether a covariate has a marginal non-linear effect. Alternatively, we can obtain some marginal estimates of pairs of two or more covariates to investigate their interactions. Furthermore, we can use the SS-ANOVA method as a diagnostic tool.<sup>37</sup>

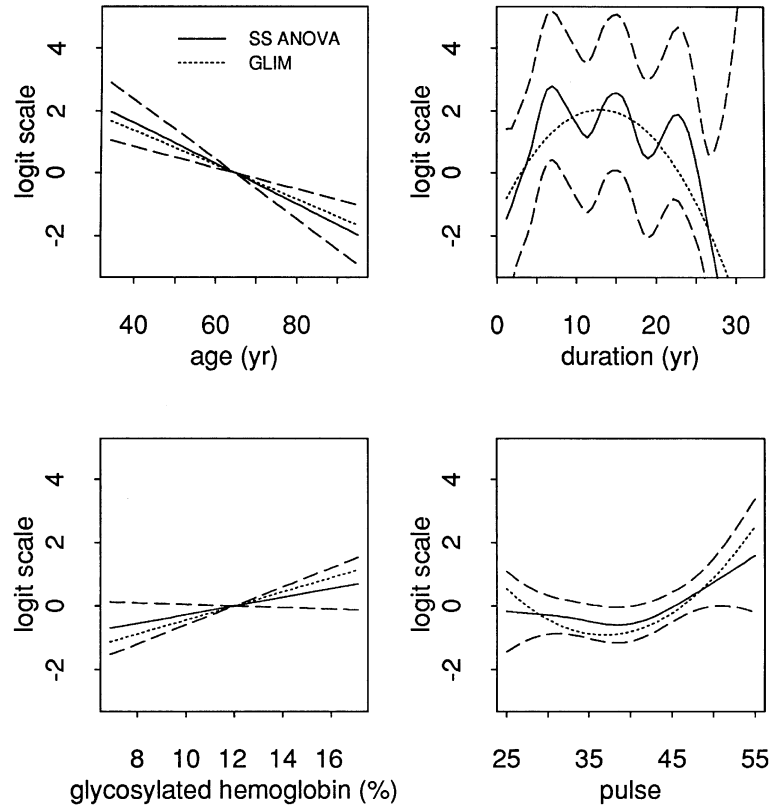


Figure 8. ID 'Incidence' analysis. Estimates of the main effects for incidence, logit scale. Dashed lines are componentwise 90 per cent Bayesian confidence intervals

After deleting participants with missing data, there were 374 participants in the older onset ID group qualified for inclusion in the analysis for 'Progression'. Klein *et al.*<sup>27</sup> found that age, duration, glycosylated haemoglobin and base-retinopathy-level are significant using a linear logistic model. They concluded that the risk of progression is higher if duration is longer and if base-retinopathy-level is lower (less severe). To investigate whether a linear logistic model is appropriate, we obtained SS-ANOVA estimates for age, duration, glycosylated haemoglobin and bmi separately. Figure 9 plots these marginal estimates and proportions within each decile on the logit scale. These plots suggest that the effects of duration and bmi may be non-linear. Of course we can only use these marginal analyses as a preliminary tool. Our goal is to build a full model with all covariates. The effects of covariates in the marginal analyses may be different from those in the full model. Using linear logistic models, we find that the effect of age is non-linear (significant up to 4th order of polynomials). The effect of bmi is borderline significant and non-linear (a 4th order polynomial has a  $p$ -value of 0.0609).

This analysis provides an opportunity to include base-retinopathy-level as a categorical variable in the model. We divided base-retinopathy-level into five categories: 10/10 as category 1; 21/< 21 and 21/21 as category 2; 31/< 31 and 31/31 as category 3; 41/< 41 and 41/41 as category 4; 51/< 51 and 51/51 as category 5. We first fit an SS-ANOVA model with main effects of age, duration, glycosylated haemoglobin, bmi and base-retinopathy-level.

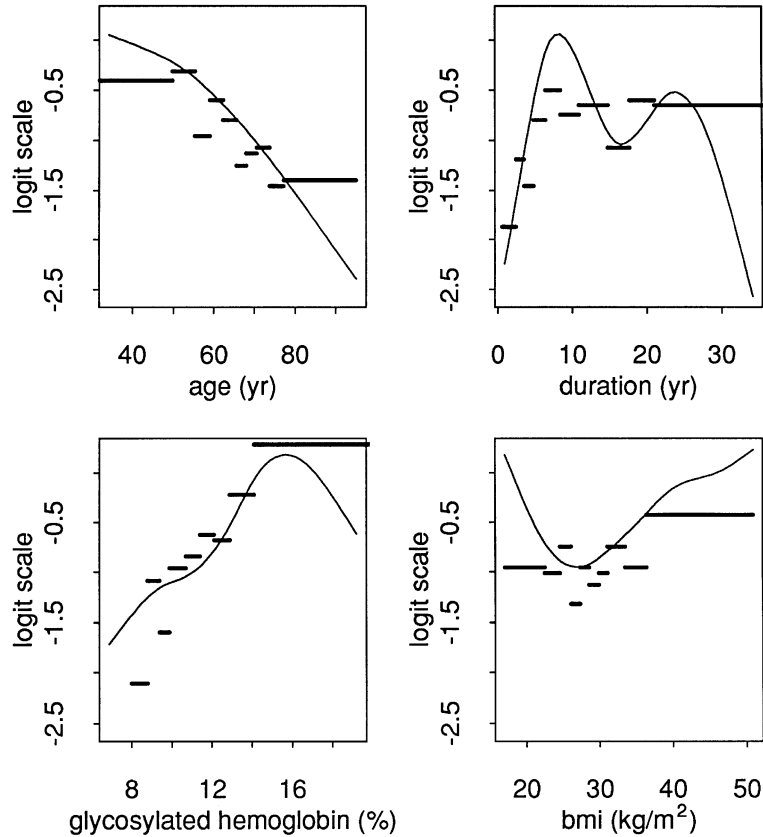


Figure 9. ID 'Progression' analysis. Separate estimates of the marginal effects. Solid lines: SS-ANOVA estimates; segments: the logit of proportions within each decile

We found that the main effect of age and glycosylated haemoglobin are linear. Finally, we fitted the model:

$f(\text{age, duration, glycosylated haemoglobin, bmi, base-retinopathy-level})$

$$\begin{aligned}
 &= \mu + a_1 \times \text{age} + f_1(\text{duration}) + a_2 \times \text{glycosylated haemoglobin} + f_2(\text{bmi}) \\
 &+ \sum_{k=2}^5 \gamma_k I_k(\text{base-retinopathy-level}).
 \end{aligned} \tag{19}$$

Figure 10 plots the estimates of the main effects and their 90 per cent Bayesian confidence intervals. The estimates of  $\gamma_2, \gamma_3, \gamma_4, \gamma_5$  and their posterior standard deviations (inside parentheses) are 0.73(0.35), -0.68(0.37), -1.05(0.40) and -0.17(0.51). Our conclusions basically agree with Klein *et al.*<sup>27</sup> with the following modifications:

1. The risk of progression of retinopathy increases with increasing duration of diabetes up to around 8 years and does not increase after that.
2. The effect of body mass index is small.
3. Baseline level has a significant effect on progression of retinopathy, but not monotonically. A subject with baseline level 10/10 has about the same risk as a subject with baseline level



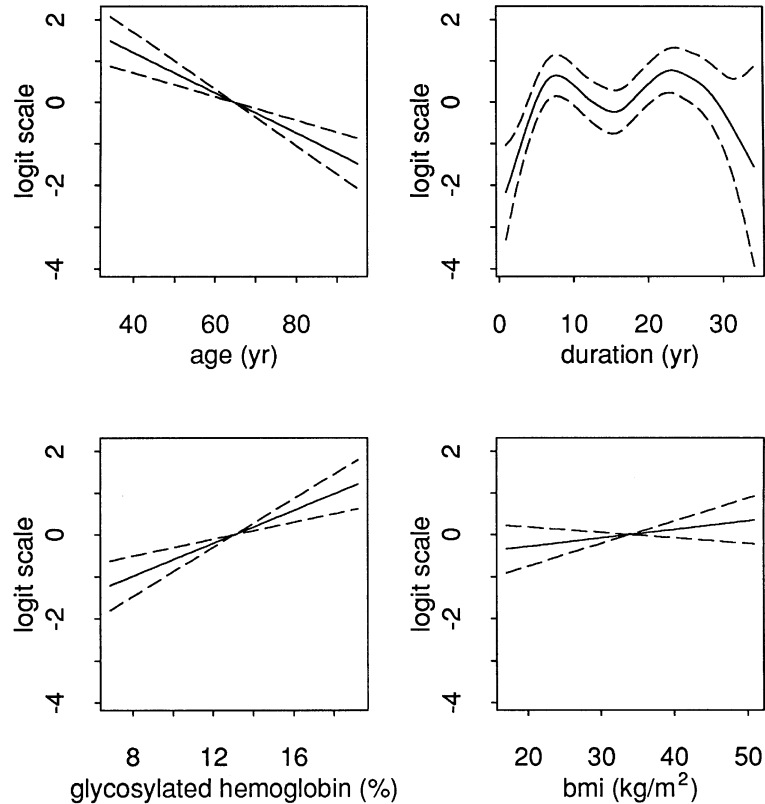


Figure 10. ID 'Progression' analysis. Estimates of the main effects for progression. Dashed lines are componentwise 90 per cent Bayesian confidence intervals

51/< 51 or 51/51. A subject with baseline level 10/10 has higher risk than a subject with baseline level 31/< 31 to 41/41. A subject with baseline level 10/10 has lower risk than a subject with baseline level 21/< 21 and 21/21. This may reflect the arbitrary division of the levels that may not be evenly spaced, or, persons with more severe retinopathy may progress more rapidly, but those who progress may less likely survive to be examined at follow-up.

## 7. CONCLUSIONS

We first summarize the new findings from our analyses, and then make some concluding remarks concerning the SS-ANOVA estimation methods.

1. Age effects are non-linear and different for different groups. For progression of the older onset ID group, the risk decreases with increasing age. Since mortality may change the population under study, especially for the older onset group, more frequent follow-up of a cohort (perhaps yearly) might provide further understanding of this relation.
2. The effect of duration is consistent in all groups. The risk generally increases with increasing duration up to a certain point and does not increase after that. This indicates that if a person with older onset diabetes has not had an event (incidence or progression of retinopathy) after some years of diabetes, the risk will not substantially increase after that, although it remains higher than that for newly diagnosed persons.

3. The effect of glycosylated haemoglobin is large, which agrees with previous studies. We find that there is a strong interaction between age and glycosylated haemoglobin for incidence in the lower onset NID group.
4. The effect of body mass index (bmi) is small.
5. Pulse rate has a moderate effect on the incidence of retinopathy in the older onset ID group. In a main effects model which also included age, duration, and glycosylated haemoglobin, the main effect for pulse rate is constant for lower pulse rates and then veers upward linearly (in the logit scale) for higher pulse rates.
6. Baseline retinopathy severity level (base-retinopathy-level) has a significant effect on progression in the older onset ID group in a main effects model that also included age, duration, glycosylated haemoglobin and body mass index, but the effect was not monotonic. Participants with baseline level  $21 / < 21$  and  $21 / 21$  had the highest risk and participants with baseline levels  $31 / < 31$  to  $41 / 41$  had the lowest risk of progression.

The SS-ANOVA models has provided us with a family of flexible penalized log-likelihood estimates, which specifically include the possible fitting of interaction terms, as well as allowing for combinations of continuous and categorical variates, which, moreover, reduce to standard GLIM models as the smoothing parameters tend to infinity. These models are well adapted to the irregular distribution of multiple predictor variables commonly found in demographic data. The data-based method for estimating smoothing parameters allows the data to suggest when the flexible or 'smooth' components of the model are unnecessary. This family of models allows the user to visualize complex relationships between variables otherwise not evident with the use of purely parametric GLIM models. For example, the analysis here leads to new questions regarding why higher glycosylated haemoglobin in the younger subjects in the NID incidence group leads to increased risk compared to higher glycosylated haemoglobin in the older subjects in this group. None of the SS-ANOVA models built for WESDR data in this paper can be easily well approximated by a parametric GLIM model. Some SS-ANOVA models (for example, incidence in the younger onset group in reference 19) can be reduced to parametric GLIM models. Therefore the same conclusions can be reached if appropriate terms are included in the parametric GLIM model. Even in these cases SS-ANOVA methods are very useful in deciding the appropriate terms.

SS-ANOVA models also provide measure of risk which the clinicians are particularly interested in.<sup>38</sup>

As with any non-parametric function estimation method, we must make various assumptions about the function estimated. The primary assumption in the models of this paper is that  $f$  is 'smooth' as measured by the integrals of  $(\partial^2 f / \partial t_\alpha^2)^2$ ,  $(\partial^2 f / \partial t_\beta^2)^2$ , and  $(\partial^4 f / \partial t_\alpha^2 \partial t_\beta^2)^2$  (which give the five terms in the penalty functional of (12)).

There is a limit to the number of smoothing parameters that we can estimate reliably with the sample sizes here. Some of the local 'wiggles' are probably a manifestation of that, particularly, in Figure 8. The Bayesian confidence intervals, however, do suggest when bumps are not 'real'. When used for exploratory purposes, one may use the smoothing parameters chosen by the automatic method here as starting guesses for 'eyeball' or subjective smoothing parameter choices.

One drawback of this family of methods (aside from the fact that one requires larger data sets than for standard purely parametric models) is the fact that the publicly available software for their use is, at the present time, not quite at the 'cookbook' level. If a sample program for the particular model of interest is unavailable, the user must write a driver that contains details of their model. Examples of drivers for three models, including the one used in WWGKK, are

packaged with GRKPACK. Interested readers may obtain drivers for the four SS-ANOVA models used in this paper from the first author. Since the GRKPACK code is based on matrix decompositions, it is relatively slow, and requires a relatively large amount of storage, placing an upper limit on the sample sizes one can analyse. Approximate numerical methods that will reduce both the time and space required for a given sample size are an area of active research at the present time. In the meantime, we believe that the present methods represent an important new approach to the analysis of demographic data sets similar to those analyzed here.

## ACKNOWLEDGEMENTS

We are grateful to the editor and two anonymous referees for helpful comments and extensive editorial suggestions. We thank Scot Moss who provided some of the data and invaluable technical support. The second author wishes to express her special appreciation to Marian Fisher who introduced her to some important issues in demographic data analysis. Yuedong Wang's research supported in part by NIH grant EY09946, P60 DK20572, P30 HD18258 and U54 HD29184, Grace Wahba's research supported in part by NIH grant EY09946 and NSF grant DMS9121003, Chong Gu's research supported by NSF grant DMS9301511, Ronald Klein's and Barbara Klein's research supported by NIH grant EY03083. References 2, 3, 19, 20, 21, 23 and 38 are available through the second author's home page URL <http://www.stat.wisc.edu/~wahba>.

## REFERENCES

1. McCullagh, P. and Nelder, J. *Generalized Linear Models*, 2nd edn, Chapman and Hall, 1989.
2. Wahba, G., Wang, Y., Gu, C., Klein, R. and Klein, B. 'Smoothing spline ANOVA for exponential families, with application to the Wisconsin Epidemiological Study of Diabetic Retinopathy', *Annals of Statistics*, **23**, 1865–1895 (1995).
3. Wang, Y. 'GRKPACK: Fitting smoothing spline analysis of variance models to data from exponential families', Technical Report 942, Department of Statistics, University of Wisconsin, Madison, WI, 1995, to appear, *Communications in Statistics – Simulation and Computation*.
4. O'Sullivan, F. The analysis of some penalized likelihood estimation schemes, PhD thesis, Department of Statistics, University of Wisconsin, Madison, WI, 1983. Technical Report 726.
5. O'Sullivan, F., Yandell, B. and Raynor, W. 'Automatic smoothing of regression functions in generalized linear models', *Journal of the American Statistical Association*, **81**, 96–103 (1986).
6. Gu, C. and Wahba, G. 'Semiparametric analysis of variance with tensor product thin plate splines', *Journal of the Royal Statistical Society, Series B*, **55**, 353–368 (1993).
7. Gu, C. and Wahba, G. 'Smoothing spline ANOVA with component-wise Bayesian "confidence intervals"', *Journal of Computational and Graphical Statistics*, **2**, 97–117 (1993).
8. Hastie, T. and Tibshirani, R. *Generalized Additive Models*, Chapman and Hall, 1990.
9. Chambers, J. and Hastie, T. *Statistical Models in S*, Wadsworth and Brooks, 1992.
10. Hastie, T. and Tibshirani, R. 'Varying-coefficient models', *Journal of the Royal Statistical Society, Series B*, **55**, 757–796 (1993).
11. Wahba, G. *Spline Models for Observational Data*, SIAM, CBMS-NSF Regional Conference Series in Applied Mathematics, v. 59, 1990.
12. Gu, C. and Wahba, G. 'Comments to "Multivariate adaptive regression splines", by J. Friedman', *Annals of Statistics*, **19**, 115–123 (1991).
13. Kimeldorf, G. and Wahba, G. 'Some results on Tchebycheffian spline functions', *Journal of Mathematical Analysis and Applications*, **33**, 82–95 (1971).
14. Wahba, G. 'Partial and interaction splines for the semiparametric estimation of functions of several variables', in Boardman, T. (ed), *Computer Science and Statistics: Proceedings of the 18th Symposium*, American Statistical Association, Washington, DC, 1986, pp. 75–80.
15. Gu, C. 'Adaptive spline smoothing in non-Gaussian regression models', *Journal of the American Statistical Association*, **85**, 801–807 (1990).
16. Gu, C. 'Cross-validating non-Gaussian data', *Journal of Computational and Graphical Statistics*, **1**, 169–179 (1992).

17. Gu, C., Bates, D. M., Chen, Z. and Wahba, G. 'The computation of GCV functions through householder tridiagonalization with application to the fitting of interaction spline models', *SIAM Journal on Matrix Analysis*, **10**, 457–480 (1989).
18. Gu, C. and Wahba, G. 'Minimizing GCV/GML scores with multiple smoothing parameters via the Newton method', *Journal of Scientific and Statistical Computing*, **12**, 383–398 (1991).
19. Wang, Y. 'Smoothing spline analysis of variance of data from exponential families', PhD thesis, Technical Report 928, University of Wisconsin-Madison, Madison, WI, 1994.
20. Wahba, G., Gu, C., Wang, Y. and Chappell R. 'Soft classification, a. k. a. risk estimation, via penalized log likelihood and smoothing spline analysis of variance', in Wolpert, D. (ed), *The Mathematics of Generalization, Santa Fe Institute Studies in the Sciences of Complexity, Proc. Vol. XX*, Addison-Wesley, Reading, MA, 1995, pp. 329–360.
21. Wang, Y., Wahba, G., Chappell, R. and Gu, C. 'Simulation studies of smoothing parameter estimates and Bayesian confidence intervals in Bernoulli SS-ANOVA models', *Communications in Statistics – Simulation and Computation*, **24**, 1037–1059 (1995).
22. Wong, W. 'Estimation of the loss of an estimate', Technical Report 356, Department of Statistics, University of Chicago, Chicago, IL, 1992.
23. Xiang, D. and Wahba, G. 'A generalized approximate cross validation for smoothing splines with non-Gaussian data', *Statistica Sinica*, **6**, 675–692 (1996).
24. Yandell, B. 'Algorithms for nonlinear generalized cross-validation', in Boardman, T. J. (ed), *Computer Science and Statistics: 18th Symposium on the Interface*, American Statistical Association, Washington, DC, 1986.
25. Wahba, G. 'Bayesian "confidence intervals" for the cross-validated smoothing spline', *Journal of the Royal Statistical Society, Series B*, **45**, 133–150 (1983).
26. Gu, C. 'Penalized likelihood regression: a Bayesian analysis', *Statistica Sinica*, **2**, 255–264 (1992).
27. Klein, R., Klein, B. E. K., Moss, S. E., Davis, M. D. and DeMets, D. L. 'Glycosylated hemoglobin predicts the incidence and progression of diabetic retinopathy', *Journal of the American Medical Association*, **260**, 2864–2871 (1988).
28. Klein, R., Klein, B. E. K., Moss, S. E., Davis, M. D. and DeMets, D. L. 'Is blood pressure a predictor of the incidence or progression of diabetic retinopathy', *Archives of Internal Medicine*, **149**, 2427–2432 (1989).
29. Klein, B. E. K., Davis, M. D., Segal, P., Long, J. A., Harris, W. A., Haug, G. A., Magli, Y. and Syrjala, S. 'Diabetic retinopathy: Assessment of severity and progression', *Ophthalmology*, **91**, 10–17 (1984).
30. Klein, R., Klein, B. E. K., Moss, S. E. and Cruickshanks, K. J. 'The relationship of hyperglycemia of long-term incidence of progression of diabetic retinopathy', *Archives of Internal Medicine*, **154**, 2169–2178 (1994).
31. Klein, R., Klein, B. E. K., Moss, S. E. and Cruickshanks, K. J. 'The Wisconsin Epidemiologic Study of Diabetic Retinopathy. XIV. Ten year incidence and progression of diabetic retinopathy', *Archives of Ophthalmology*, **112**, 1217–1228 (1994).
32. Klein, R., Klein, B. E. K., Moss, S. E., Davis, M. D. and DeMets, D. L. 'The Wisconsin Epidemiologic Study of Diabetic Retinopathy. II. Prevalence and risk of diabetic retinopathy when age at diagnosis is less than 30 years', *Archives of Ophthalmology*, **102**, 520–526 (1984).
33. Klein, R., Klein, B. E. K., Moss, S. E., Davis, M. D. and DeMets, D. L. 'The Wisconsin Epidemiologic Study of Diabetic Retinopathy. III. Prevalence and risk of diabetic retinopathy when age at diagnosis is 30 or more years', *Archives of Ophthalmology*, **102**, 527–532 (1984).
34. Klein, R., Klein, B. E. K., Moss, S. E., Davis, M. D. and DeMets, D. L. 'The Wisconsin Epidemiologic Study of Diabetic Retinopathy. IX. Four year incidence and progression of diabetic retinopathy when age at diagnosis is less than 30 years', *Archives of Ophthalmology*, **107**, 237–243 (1989).
35. Klein, R., Klein, B. E. K., Moss, S. E., Davis, M. D. and DeMets, D. L. 'The Wisconsin Epidemiologic Study of Diabetic Retinopathy. X. Four incidence and progression of diabetic retinopathy when age at diagnosis is 30 or more years', *Archives of Ophthalmology*, **107**, 244–249 (1989).
36. Moss, S. E., Klein, R., Klein, B. E. K., Meuer, S. M. 'The association of glycemia and cause-specific mortality in a diabetic population', *Archives of Internal Medicine*, **154**, 2473–2479 (1994).
37. Fowlkes, E. B. 'Some diagnostics for binary logistic regression via smoothing', *Biometrika*, **74**, 503–515 (1987).
38. Wang, Y. 'Odds ratio estimation in Bernoulli smoothing spline ANOVA model', Technical Report 946, University of Wisconsin-Madison, Department of Statistics, 1996, to appear, *The Statistician*.