

**Smoothing Spline Mixed-Effects Density Models
for Clustered Data**

Chi-Yang Chiu, Anna Liu and Yuedong Wang

University of Tennessee Health Science Center,

University of Massachusetts at Amherst, and

University of California at Santa Barbara

Supplementary Material

**S1 An example on construction of an NMECDM based
on an SS ANOVA decomposition**

We assume $\mathcal{Y} = \mathbb{R}$ and model g as a function of y using the thin-plate spline model space $\mathcal{H}_y = W_2^3(\mathbb{R}) \ominus \{1\}$. Suppose that we want to model g as a function of x using an RKHS $\mathcal{H}_x = \mathcal{H}_{0x} \oplus \mathcal{H}_{1x}$ where \mathcal{H}_{0x} is a finite dimensional space collecting functions which are not penalized. Let P_y and P_x be the projection operators onto \mathcal{H}_{0y} and \mathcal{H}_{0x} respectively. Let

$P_\omega g = \int_\Omega g(y, x, \omega) dP$. We have the following SS ANOVA decomposition

$$\begin{aligned}
 g &= [P_y + (I - P_y)][P_x + (I - P_x)][P_\omega + (I - P_\omega)]g \\
 &\stackrel{\Delta}{=} g_{ppf}(y, x) + g_{psf}(y, x) + g_{spf}(y, x) + g_{ssf}(y, x) \\
 &\quad + g_{ppr}(y, x, \omega) + g_{psr}(y, x, \omega) + g_{spr}(y, x, \omega) + g_{ssr}(y, x, \omega) \quad (\text{S1.1})
 \end{aligned}$$

where g_{ppf} is the parametric fixed effect, g_{psf} , g_{spf} and g_{ssf} are nonparametric fixed effects, g_{ppr} is the parametric random effect, and g_{psr} , g_{spr} and g_{ssr} are nonparametric random effects. Comparing to the NMECDM (2.5), we have $\mathcal{H}_\eta = (W_2^3(\mathbb{R}) \ominus \{1\}) \otimes \mathcal{H}_x$, $\eta(y, x) = g_{ppf}(y, x) + g_{psf}(y, x) + g_{spf}(y, x) + g_{ssf}(y, x)$ and $b_i(y) = g_{ppr}(y, x, \omega_i) + g_{psr}(y, x, \omega_i) + g_{spr}(y, x, \omega_i) + g_{ssr}(y, x, \omega_i)$.

Rather than the conditional density, regression models focus on the conditional expectation $\mu(x, \omega) = E(Y|X = x, \omega)$. Suppose that we want to model μ as a function of x using the RKHS $\mathcal{H}_x = \mathcal{H}_{0x} \oplus \mathcal{H}_{1x}$. Applying a similar SS ANOVA decomposition to the random function $\mu(x, \omega)$ we have the following SS ANOVA mixed effects model (Wang, 1998)

$$Y_{ij} = \mu_{pf}(X_{ij}) + \mu_{sf}(X_{ij}) + \mu_{pr}(X_{ij}, \omega_i) + \mu_{sr}(X_{ij}, \omega_i) + \epsilon_{ij}, \quad (\text{S1.2})$$

where μ_{pf} and μ_{sf} are parametric and nonparametric fixed effects, μ_{pr} and μ_{sr} are parametric and nonparametric random effects, $\epsilon_{ij} \stackrel{iid}{\sim} N(0, \sigma^2)$, and random effects and ϵ_{ij} are mutually independent. Then, up to a con-

stant independent of y , the logistic density of Y_{ij} conditional on $X_{ij} = x$ and cluster ω has the form $\{(-y^2/2 + \mu_{pf}y) + \mu_{sf}(x)y + \mu_{pr}(x, \omega)y + \mu_{sr}(x, \omega)y\}/\sigma^2$. Comparing to the SS ANOVA decomposition (S1.1), it is obvious that the SS ANOVA mixed effects model (S1.2) is a special case with $g_{ppf}(y, x) \cong (-y^2/2 + \mu_{pf}(x)y)/\sigma^2$, $g_{psf}(y, x) \cong \mu_{sf}(x)y/\sigma^2$, $g_{ppr}(y, x, \omega) \cong \mu_{pr}(x, \omega)y/\sigma^2$, $g_{psr}(y, x, \omega) \cong \mu_{sr}(x, \omega)y/\sigma^2$, and $g_{spf}(y, x, \omega) = g_{ssf}(y, x, \omega) = g_{spr}(y, x, \omega) = g_{ssr}(y, x, \omega) = 0$.

S2 Derivatives of log-likelihood

Note that

$$\log f_{\mathbf{Y}_i | \mathbf{B}_i}(\mathbf{Y}_i) \cong \sum_{j=1}^{n_i} \left\{ \eta(Y_{ij}, X_{ij}) + b_i(Y_{ij}, X_{ij}) - \log \int_{\mathcal{Y}} \exp\{\eta(y, X_{ij}) + b_i(y, X_{ij})\} dy \right\},$$

where

$$\eta(y, x) = \sum_{\nu=1}^p d_{\nu} \phi_{\nu}(y, x) + \sum_{l=1}^L c_l R_1^*(\mathbf{U}_l, (y, x)).$$

Hence

$$\begin{aligned} & \frac{\partial \log f_{\mathbf{Y}_i | \mathbf{B}_i}(\mathbf{Y}_i)}{\partial c_l} \\ &= \sum_{j=1}^{n_i} \left\{ R_1^*(\mathbf{U}_l, (Y_{ij}, X_{ij})) - \frac{\int_{\mathcal{Y}} R_1^*(\mathbf{U}_l, (y, X_{ij})) \exp\{\eta(y, X_{ij}) + b_i(y, X_{ij})\} dy}{\int_{\mathcal{Y}} \exp\{\eta(y, X_{ij}) + b_i(y, X_{ij})\} dy} \right\} \\ &= \sum_{j=1}^{n_i} \left\{ R_1^*(\mathbf{U}_l, (Y_{ij}, X_{ij})) - E_{\mathbf{Y}_i | \mathbf{B}_i} R_1^*(\mathbf{U}_l, (Y, X_{ij})) \right\}, \end{aligned}$$

and

$$\begin{aligned}
 & \frac{\partial^2 \log f_{\mathbf{Y}_i | \mathbf{B}_i}(\mathbf{Y}_i)}{\partial c_l \partial c_k} \\
 = & - \sum_{j=1}^{n_i} \left\{ \frac{\int_{\mathcal{Y}} R_1^*(\mathbf{U}_l, (y, X_{ij})) R_1^*(\mathbf{U}_k, (y, X_{ij})) \exp\{\eta(y, X_{ij}) + b_i(y, X_{ij})\} dy}{\int_{\mathcal{Y}} \exp\{\eta(y, X_{ij}) + b_i(y, X_{ij})\} dy} \right. \\
 & - \left(\frac{\int_{\mathcal{Y}} R_1^*(\mathbf{U}_l, (y, X_{ij})) \exp\{\eta(y, X_{ij}) + b_i(y, X_{ij})\} dy}{\int_{\mathcal{Y}} \exp\{\eta(y, X_{ij}) + b_i(y, X_{ij})\} dy} \right) \\
 & \left. \left(\frac{\int_{\mathcal{Y}} R_1^*(\mathbf{U}_k, (y, X_{ij})) \exp\{\eta(y, X_{ij}) + b_i(y, X_{ij})\} dy}{\int_{\mathcal{Y}} \exp\{\eta(y, X_{ij}) + b_i(y, X_{ij})\} dy} \right) \right\} \\
 = & \sum_{j=1}^{n_i} \text{Cov}_{\mathcal{Y} | \mathbf{B}_i} (R_1^*(\mathbf{U}_k, (Y, X_{ij})), R_1^*(\mathbf{U}_l, (Y, X_{ij}))).
 \end{aligned}$$

Other first and second derivatives can be calculated similarly.

S3 Markov Chain Monte Carlo

We need to generate MCMC samples for the computation of conditional expectations with respect to $\mathbf{B}_i | \mathbf{Y}_i$. We first divide the domain \mathcal{Y} into a finite number of disjoint subsets $\mathcal{Y} = \cup_{k=1}^K \mathcal{Y}_k$ and select points $\tilde{y}_k \in \mathcal{Y}_k$ for $k = 1, \dots, K$. We then apply the Metropolis-Hastings (MH) procedure to generate MCMC samples for random vectors $\tilde{\mathbf{B}}_{ij} = (b_i(\tilde{y}_1, X_{ij}), \dots, b_i(\tilde{y}_K, X_{ij}))$ conditional on \mathbf{Y}_i for $i = 1, \dots, m$ and $j = 1, \dots, n_i$. Finally we approximate the conditional expectations with respect to $\mathbf{B}_i | \mathbf{Y}_i$ using these MCMC samples. Details of the MH procedure can be found in Gelman et al. (2003). We use a multivariate normal centered at the current MCMC

sample as the proposal distribution with a scaled covariance matrix such that the acceptance rate is near 23% as suggested by Gelman et al. (2003, Ch11). Convergence and mixing are assessed by visualization tools such as trace plot and autocorrelation.

S4 Quadratic Approximation and Cross-Validation

The log marginal likelihood for cluster i is $l_i = \log E_{\mathbf{B}_i} f_{\mathbf{Y}_i | \mathbf{B}_i}(\mathbf{Y}_i)$ where $f_{\mathbf{Y}_i | \mathbf{B}_i}(\mathbf{Y}_i)$ is given in (3.2). For $f, g \in \mathcal{H}_\eta$ and $\alpha \in \mathbb{R}$, consider

$$L_{f,g}(\alpha) = \log E_{\mathbf{B}_i} \frac{\exp\{\sum_{j=1}^{n_i} [f(Y_{ij}, X_{ij}) + \alpha g(Y_{ij}, X_{ij}) + b_i(Y_{ij}, X_{ij})]\}}{\int_{\mathcal{Y}} \exp\{f(y, X_{ij}) + \alpha g(y, X_{ij}) + b_i(y, X_{ij})\} dy} \quad (\text{S4.1})$$

as a function of α . It can be shown that the derivatives of $L_{f,g}(\alpha)$ with respect to α evaluated at zero are given as follows:

$$\begin{aligned} L'_{f,g}(0) &= \sum_{j=1}^{n_i} g(Y_{ij}, X_{ij}) - \sum_{j=1}^{n_i} E_{\mathbf{B}_i | \mathbf{Y}_i}^f E_{\mathbf{Y} | \mathbf{B}_i}^f g(Y, X_{ij}), \\ L''_{f,g}(0) &= \sum_{j=1}^{n_i} E_{\mathbf{B}_i | \mathbf{Y}_i}^f V_{\mathbf{Y} | \mathbf{B}_i}^f g(Y, X_{ij}) + V_{\mathbf{B}_i | \mathbf{Y}_i}^f \left\{ \sum_{j=1}^{n_i} E_{\mathbf{Y} | \mathbf{B}_i}^f g(Y, X_{ij}) \right\}, \end{aligned}$$

where $E_{\cdot | \cdot}^f$ and $V_{\cdot | \cdot}^f$ respectively represent conditional expectation and variance under model (2.5) with $\eta(y, x) = f(y, x)$.

For any fixed $\tilde{\eta}$, setting $f = \tilde{\eta}$ and $g = \eta - \tilde{\eta}$, we have the quadratic

approximation of the log-likelihood l_i at $\tilde{\eta}$

$$\begin{aligned}
 l_i &= L_{f,g}(1) \\
 &\approx L_{f,g}(0) + L'_{f,g}(0) + \frac{1}{2}L''_{f,g}(0) \\
 &= \log \mathbb{E}_{\mathbf{B}_i} \frac{\exp\{\sum_{j=1}^{n_i} [\tilde{\eta}(Y_{ij}, X_{ij}) + b_i(Y_{ij}, X_{ij})]\}}{\int_{\mathcal{Y}} \exp\{\tilde{\eta}(y, X_{ij}) + b_i(y, X_{ij})\} dy} + \\
 &\quad \sum_{j=1}^{n_i} \{\eta(Y_{ij}, X_{ij}) - \tilde{\eta}(Y_{ij}, X_{ij})\} \\
 &\quad - \sum_{j=1}^{n_i} \mathbb{E}_{\mathbf{B}_i | \mathbf{Y}_i}^{\tilde{\eta}} \mathbb{E}_{Y | \mathbf{B}_i}^{\tilde{\eta}} \{\eta(Y, X_{ij}) - \tilde{\eta}(Y, X_{ij})\} \\
 &\quad + \frac{1}{2} \sum_{j=1}^{n_i} \mathbb{E}_{\mathbf{B}_i | \mathbf{Y}_i}^{\tilde{\eta}} \mathbb{V}_{Y | \mathbf{B}_i}^{\tilde{\eta}} \{\eta(Y, X_{ij}) - \tilde{\eta}(Y, X_{ij})\} \\
 &\quad + \frac{1}{2} \mathbb{V}_{\mathbf{B}_i | \mathbf{Y}_i}^{\tilde{\eta}} \left\{ \sum_{j=1}^{n_i} \mathbb{E}_{Y | \mathbf{B}_i}^{\tilde{\eta}} (\eta(Y, X_{ij}) - \tilde{\eta}(Y, X_{ij})) \right\}.
 \end{aligned}$$

Dropping terms do not involve η and the last term for computational stability, and summing log-likelihoods over all clusters, we have the quadratic approximation to log marginal likelihood at $\tilde{\eta}$

$$\begin{aligned}
 \tilde{l}(\boldsymbol{\zeta}, \eta) &= \sum_{i=1}^m \sum_{j=1}^{n_i} \left\{ \eta(Y_{ij}, X_{ij}) - \mathbb{E}_{\mathbf{B}_i | \mathbf{Y}_i}^{\tilde{\eta}} \mathbb{E}_{Y | \mathbf{B}_i}^{\tilde{\eta}} \eta(Y, X_{ij}) \right. \\
 &\quad \left. + \frac{1}{2} \mathbb{E}_{\mathbf{B}_i | \mathbf{Y}_i}^{\tilde{\eta}} \mathbb{V}_{Y | \mathbf{B}_i}^{\tilde{\eta}} (\eta(Y, X_{ij}) - \tilde{\eta}(Y, X_{ij})) \right\}.
 \end{aligned}$$

Consider the approximated penalized likelihood

$$-\frac{1}{N} \tilde{l}(\boldsymbol{\zeta}, \eta) + \frac{1}{2} \sum_{j=1}^q \lambda_j \|P_j \eta\|^2, \quad (\text{S4.2})$$

and denote $\eta_{\boldsymbol{\lambda}}^{[i,j]}$ as the solution to (S4.2) with the j th observation from cluster i begin removed. Let $\boldsymbol{\psi} = (\phi_1, \dots, \phi_p, \xi_1, \dots, \xi_L)^T$, $\check{R} = (\boldsymbol{\psi}(\mathbf{Z}_{11}), \dots,$

$\boldsymbol{\psi}(\mathbf{Z}_{mnm}))^T$ where $\mathbf{Z}_{ij} = (Y_{ij}, X_{ij})$, $P_1^\perp = I_N - \mathbf{1}_N \mathbf{1}_N^T / N$ where I_N is an $N \times N$ identity matrix and $\mathbf{1}_N$ is an N -vector of all ones, and $\Pi = E_{\mathbf{B}|\mathbf{Y}} I(\tilde{\mathbf{c}}, \tilde{\mathbf{d}}, \mathbf{B})$ where I is a proxy of the Hessian matrix defined in Section 3.2, and $\tilde{\mathbf{c}}$ and $\tilde{\mathbf{d}}$ are estimates at the current iteration. Following similar arguments as in Gu (2013), it can be shown that

$$\begin{aligned} \eta_{\boldsymbol{\lambda}}^{[i,j]}(Y_{ij}, X_{ij}) &= \eta_{\boldsymbol{\lambda}}(Y_{ij}, X_{ij}) - \frac{1}{N-1} \{\boldsymbol{\psi}(\mathbf{Z}_{ij})\}^T \Pi^{-1} \{\boldsymbol{\psi}(\mathbf{Z}_{ij}) \\ &\quad - N^{-1} \check{R}^T \mathbf{1}\}. \end{aligned}$$

Noting that $\check{R}^T \mathbf{1} = \sum_{i=1}^m \sum_{j=1}^{n_i} \boldsymbol{\psi}(\mathbf{Z}_{ij})$, we have

$$\frac{1}{N} \sum_{i=1}^m \sum_{j=1}^{n_i} \eta_{\boldsymbol{\lambda}}^{[i,j]}(Y_{ij}, X_{ij}) = \frac{1}{N} \sum_{i=1}^m \sum_{j=1}^{n_i} \eta_{\boldsymbol{\lambda}}(Y_{ij}, X_{ij}) - \frac{\text{tr}(P_1^\perp \check{R} \Pi^{-1} \check{R}^T P_1^\perp)}{N(N-1)}.$$

The cross-validation score (4.4) is derived by plugging in the above into (4.3).

S5 Figures and A Table

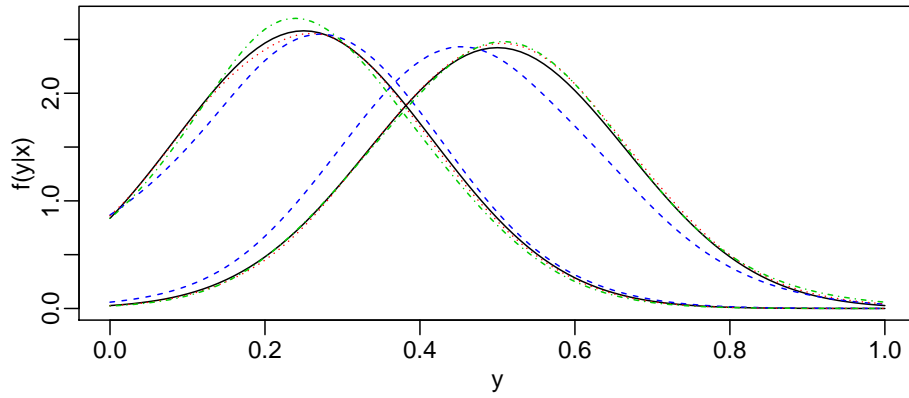


Figure A.1: True and estimates of the population conditional density function when x is discrete and $m = 200$. The solid black curves are the true population conditional density functions of the two groups, the dashed blue curves are the estimates with the largest AKL loss, the dotted red curves are the estimates with the smallest AKL loss, and the dash-dot green curves are the estimates with the median AKL loss.

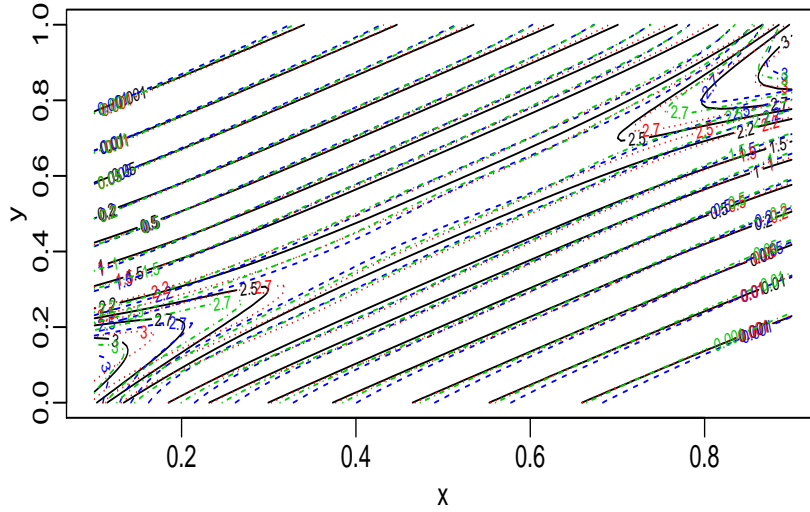


Figure A.2: Contour plots of the true population conditional density function (solid black lines) and three estimates corresponding to the largest AKL loss (dashed blue lines), median AKL loss (dash-dot green lines), and the smallest AKL loss (dotted red lines) when x is continuous and $m = 200$.

Parameter	Subset 1	Subset 2	Subset 3	Subset 4	Subset 5	Subset 6	Subset 7	Subset 8
σ_1^2	0.75	0.31	0.82	0.71	0.83	0.84	0.74	0.71
σ_2^2	142.27	192.23	144.71	142.38	144.11	144.97	143.37	143.21

Table A.1: Estimates of the variance components from the eight subsets of the Hb measurements.