

## برآورد ماکسیمم درست‌نمایی پارامترها در مدل خطی تابعی تعمیم یافته

فاطمه پای، پرویز ملک زاده، فاطمه حسینی<sup>۱</sup>

تاریخ دریافت: ۹۷/۱۰/۱۸

تاریخ پذیرش: ۹۹/۳/۸

چکیده:

گاهی در عمل داده‌ها به صورت تابعی از یک متغیر دیگر هستند که به این نوع داده‌ها، داده‌های تابعی گفته می‌شود. اگر متغیر پاسخ اسکالر و به صورت رسته‌ای یا گسسته باشد و متغیرهای کمکی به صورت تابعی، آنگاه برای تحلیل این نوع داده‌ها از مدل خطی تابعی تعمیم یافته استفاده می‌شود. در این مقاله یک مدل بریده شده خطی تابعی تعمیم یافته بررسی و برای به دست آوردن برآورد پارامترهای مدل از یک رهیافت ماکسیمم درست‌نمایی استفاده می‌شود. در نهایت در یک مطالعه شبیه‌سازی و دو مثال کاربردی مدل و روش‌های ارائه شده پیاده‌سازی می‌شوند.

**واژه‌های کلیدی:** عملگر کوواریانس، توابع ویژه، رگرسیون تابعی، مدل خطی تابعی تعمیم یافته، بسط کارهونن لونه‌و.

نمایی فرض می‌شود و مدل به صورت  $Y_i = g(\alpha + \int \beta(t)X_i(t)dt) + e_i$  بیان می‌شود که  $\alpha$  عرض از مبدأ و  $\beta(t)$  توابع پارامتری رگرسیونی نامیده می‌شوند. مدل‌های خطی تابعی توسط افرادی مثل رمسی و سیلورمن [۱۳]، فاراوی [۸]، کاردوت و همکاران [۲] و... مورد مطالعه قرار گرفته‌اند. در این مقاله مدل خطی تابعی تعمیم یافته‌ای مورد بررسی قرار گرفته است که متغیر پاسخ اسکالر و به صورت رسته‌ای یا گسسته و متغیرهای کمکی به شکل تابعی باشند. برای مدل بندی توابع و برآورد پارامترها از یک مدل بریده شده خطی تابعی تعمیم یافته استفاده و از رهیافت شبه درست‌نمایی به بررسی این مدل پرداخته شده است. ساختار مقاله به این صورت است که در بخش دوم به معرفی مدل‌های خطی تابعی تعمیم یافته و یک مدل خاص خطی تابعی تعمیم یافته پرداخته می‌شود. در بخش سوم برآورد پارامترهای مدل با رهیافت شبه درست‌نمایی بررسی و با استفاده از یک رهیافت مجانبی آزمون فرضیه و بازه اطمینان تقریبی هم‌زمان برای پارامترهای تابعی رگرسیونی به دست آورده می‌شود. در بخش چهارم در خصوص انتخاب تعداد رتبه در مدل بریده شده بحث و در نهایت در بخش پنجم بر روی چند مثال مدل و روش برآورد پیاده‌سازی می‌شود.

## ۱ مقدمه

گاهی در عمل متغیر کمی به صورت گسسته یا رسته‌ای است که مدل‌های خطی برای تحلیل این نوع مسائل معتبر نیست و از مدل‌های خطی تعمیم یافته<sup>۲</sup> استفاده می‌شود. اولین بار نلدر و ودربرن [۱۲] و به‌طور مفصل تر مک کلاو و نلدر [۹] مدل‌های خطی تعمیم یافته را معرفی نمودند و این مدل‌ها را برای مدل بندی متغیرهای پاسخ رسته‌ای و گسسته پیشنهاد دادند. چون اکثر توزیع‌های گسسته و پیوسته عضو خانواده نمایی هستند در مدل‌های خطی تعمیم یافته متغیر پاسخ متعلق به خانواده نمایی فرض می‌شود. گاهی در عمل با داده‌هایی سروکار داریم که به‌طور ذاتی به صورت تابعی از یک متغیر دیگر مانند زمان،... هستند که به این نوع داده‌ها، داده‌های تابعی می‌گویند. ممکن است در عمل با مطالعاتی مواجه شویم که متغیر پاسخ گسسته و متغیرهای کمکی به صورت تابعی هستند که در این صورت از مدل‌های خطی تابعی تعمیم یافته<sup>۳</sup> که کلاس بزرگ‌تری از مدل خطی تابعی و شامل آن است استفاده می‌شود. برای تعریف این مدل معمولاً فرض می‌کنند متغیر وابسته  $Y_i, i = 1, \dots, n$  یک متغیر تصادفی با مقادیر اسکالر است که پیوسته یا گسسته می‌باشد و متغیر کمکی به صورت یک منحنی تصادفی  $\{X_i(t), t \in \mathcal{T}\}, i = 1, \dots, n$  متناظر با یک فرایند تصادفی با میانگین صفر و روی فاصله  $\mathcal{T}$  است. مشابه مدل خطی تعمیم یافته در مدل خطی تابعی تعمیم یافته توزیع متغیر پاسخ متعلق به خانواده

<sup>۱</sup> دانشگاه سمنان، گروه آمار [fatemeh.hoseini@semnan.ac.ir](mailto:fatemeh.hoseini@semnan.ac.ir)

<sup>۲</sup> Generalized Linear Model

<sup>۳</sup> Generalized Functional Linear Model

## ۲ مدل خطی تابعی تعمیم یافته

همچنین  $E(e) = 0$ . قرار دهید  $j = 1, 2, \dots$  یک پایه یکامتعامل از فضای تابعی هیلبرت  $L^2(d\omega)$  که

$$\int_{\tau} \rho_j(t) \rho_k(t) d\omega(t) = \delta_{jk},$$

آنگاه فرآیند پیشگویی  $X(t)$  و تابع پارامتر  $\beta(t)$  با استفاده از قضیه کارهونن-لونه می توانند به صورت

$$\beta(t) = \sum_{j=1}^{\infty} \beta_j \rho_j(t)$$

$$X(t) = \sum_{j=1}^{\infty} \varepsilon_j \rho_j(t),$$

بسط داده شوند، که در آن  $E(\varepsilon_j) = 0$

$$\varepsilon_j = \int X(t) \rho_j(t) d\omega(t),$$

$$\beta_j = \int \beta(t) \rho_j(t) d\omega(t), \quad \sum \beta_j^2 < \infty.$$

با در نظر گرفتن  $\sigma_j^2 = E(\varepsilon_j^2)$

$$\sum \sigma_j^2 = \int E(X^2(t)) d\omega(t) < \infty.$$

اکنون با توجه به توابع یکامتعامل  $\rho_j$  می توان نتیجه گرفت که

$$\int \beta(t) X(t) d\omega(t) = \sum_{j=1}^{\infty} \beta_j \varepsilon_j.$$

معمولاً از خطاهای استاندارد شده که به صورت

$$e' = e\sigma(\mu) = e\bar{\sigma}(\eta)$$

هستند استفاده می شود، به طوری که  $E(e' | X) = 0$ ،  $E(e') = 0$  و  $E(e'^2) = 1$

و فرض می کنیم  $E(e'^2) = \mu^2 < \infty$ . در مدل (۱) ساختار خطاها  $(e_i)$  لازم نیست که حتماً متعلق به خانواده نمایی باشد.

### ۱.۲ مدل بریده شده خطی تابعی تعمیم یافته

در این مدل از ایده برش<sup>۴</sup> استفاده می شود. با تعریف

$$V_p = \sum_{j=p+1}^{\infty} \beta_j \varepsilon_j$$

$$U_p = \alpha + \sum_{j=1}^p \beta_j \varepsilon_j$$

آنگاه می توان نوشت

$$E(Y | X(t), t \in \tau) = g(\alpha + \sum_{j=1}^{\infty} \beta_j \varepsilon_j) = g(U_p + V_p).$$

با شرطی کردن روی  $p$  مؤلفه اول و به کار بردن تابع توزیع شرطی  $dF_{V_p|U_p}$  منجر به تابع پیوند بریده  $g_p$  به صورت

$$E(Y | U_p) = g_p(U_p) = E[g(U_p + V_p) | U_p]$$

$$= \int g(U_p + s) dF_{V_p|U_p}(s)$$

اگر هدف بررسی ارتباط متغیر پاسخ رسته ای با برخی متغیرهای کمکی تابعی است، از مدل خطی تابعی تعمیم یافته که کلاس بزرگ تری از مدل خطی تابعی و شامل مدل خطی تابعی است، می توان استفاده نمود. فرض کنید مشاهده برای  $i$  امین واحد آزمایش به صورت  $(\{X_i(t), t \in \tau\}, Y_i)$ ،  $i = 1, \dots, n$  باشد که مستقل و هم توزیع هستند. متغیر کمکی  $t \in \tau$  و  $X(t)$  به صورت یک منحنی تصادفی است و متناظر با یک فرآیند تصادفی با میانگین صفر و روی فاصله  $\tau$  است. متغیر وابسته یا پاسخ  $Y$  یک متغیر تصادفی با مقادیر اسکالر است که می تواند پیوسته یا گسسته باشد. برای مثال حالتی خاص از رگرسیون تابعی دو جمله ای، یعنی  $Y \in \{0, 1\}$  داشته باشیم. فرض کنید  $g(\cdot)$  یک تابع پیوند یکنوا مشتق پذیر معلوم است،  $\beta(\cdot)$  پارامترهای رگرسیونی تابعی و  $\sigma^2(\cdot)$  واریانس باشند در این صورت با تعیین یک اندازه مثل  $d\omega$  روی  $\tau$ ، پیشگویی خطی مدل به صورت

$$\eta = \alpha + \int \beta(t) X(t) d\omega(t)$$

تعریف می شود که

$$E(Y | X(t), t \in \tau) = \mu, \quad \mu = g(\eta)$$

$$\text{Var}(Y | X(t), t \in \tau) = \sigma^2(\mu) = \bar{\sigma}^2(\eta)$$

هستند. مشابه مدل خطی تعمیم یافته در مدل خطی تعمیم یافته تابع توزیعی  $Y$  متعلق به خانواده نمایی فرض می شود. مدل خطی تعمیم یافته تابعی معمولاً به صورت

$$Y_i = g(\alpha + \int \beta(t) X_i(t) d\omega(t)) + e_i, \quad i = 1, \dots, n \quad (1)$$

بیان می شود که در آن فرض می شود

$$E(e | X(t), t \in \tau) = 0$$

$$\text{Var}(e | X(t), t \in \tau) = \sigma^2(\mu) = \bar{\sigma}^2(\eta)$$

و  $\alpha$  یک ثابت عرض از مبدأ است و برای هر  $t \in \tau$   $E(X(t)) = 0$  و خطاها  $(e_i)$  مستقل هستند و انتگرال با توجه به اندازه  $d\omega(t)$  برای توازن وزن نامنفی  $v(\cdot)$  به طوری که اگر  $t \in \tau$  و  $v(t) > 0$  و اگر  $t \notin \tau$  و  $v(t) = 0$   $d\omega(t) = v(t)dt$  و  $v(t) = 1_{\{t \in \tau\}}$  معمولاً  $d\omega(t) = dt$  است. فرض کنید  $\sigma^2 = E\{\bar{\sigma}^2(\eta)\}$  آنگاه

$$\text{Var}(e) = \text{Var}\{E(e | X(t), t \in \tau)\} + E(\text{Var}(e | X(t), t \in \tau))$$

$$= E\{\bar{\sigma}^2(\eta)\} = \sigma^2$$

<sup>4</sup>Truncation

حل شود. با فرض معلوم بودن  $\sigma^2$  و  $g$  و با در نظر گرفتن

$$\varepsilon^{(i)} = (\varepsilon_0^{(i)}, \dots, \varepsilon_p^{(i)})'$$

$$\eta_i = \sum_{j=0}^p \beta_j \varepsilon_j^{(i)}, \mu_i = g(\eta_i), i = 1, \dots, n$$

تابع امتیاز بردار-مقدار به صورت

$$U(\beta) = \sum_{i=1}^n (Y_i - \mu_i) g'(\eta_i) \varepsilon^{(i)} / \sigma^2(\mu_i) \quad (5)$$

تعریف می شود، (مولر و استدمولر، [۱۱]). جواب معادله امتیاز (۵) به صورت

$$\hat{\beta}' = (\hat{\beta}_0, \dots, \hat{\beta}_p), \hat{\alpha} = \hat{\beta}_0 \quad (6)$$

فرض می شوند. برای حل راحت تر معادله امتیاز (۵) ماتریس های

$$D = D_{n,p} = (g'(\eta_i) \varepsilon_k^{(i)} / \sigma(\mu_i)), \quad 1 \leq i \leq n, \quad 0 \leq k \leq p,$$

$$V = V_{n,p} = \text{diag}(\sigma^2(\mu_1), \dots, \sigma^2(\mu_n)), \quad 1 \leq i \leq n,$$

تعریف می شوند. و با به کار بردن  $\mu, \varepsilon, \eta$  و تعریف

$$\begin{aligned} \Gamma &= \Gamma_p = (\gamma_{kl})_{0 \leq k, l \leq p}, \\ \gamma_{kl} &= E\left(\frac{g''(\eta)}{\sigma^2(\mu)} \varepsilon_l \varepsilon_k\right), \\ \Gamma^{-1} &= (\xi_{kl})_{0 \leq k, l \leq p} \end{aligned} \quad (7)$$

می توان نوشت

$$\Gamma = \frac{1}{n} E(D'D)$$

که یک ماتریس متقارن معین مثبت است و معکوس آن وجود دارد، که در غیر این صورت

$$E\left(\left(\sum_{k=0}^p \alpha_k \varepsilon_k g'(\eta) / \sigma(\mu)\right)^2\right) = 0$$

برای ثابت های غیر صفر  $\alpha_0, \dots, \alpha_p$  برقرار نمی شود. اکنون با در نظر گرفتن بردارهای  $Y' = (Y_1, \dots, Y_n)$ ،  $\mu' = (\mu_1, \dots, \mu_n)$  معادله  $U(\beta) = 0$  می توان به صورت ماتریس

$$D'V^{-1}(Y - \mu) = 0$$

نوشت. این معادله معمولاً با روش های تکرار حل می شود و  $\beta$ ها برآورد می شوند. تحت فرض  $\frac{1}{n}E(D'D) = \Gamma_p$  یک ماتریس معین مثبت برای هر  $p$  است، آنگاه برای هر  $p$  ثابت، یک راه حل یکتا وجود خواهد داشت. در روش بالا  $g(\cdot), \sigma^2(\cdot)$  هر دو معلوم فرض می شوند، اما می توان حالتی را فرض کرد که هر دو نامعلوم هستند که با در نظر اینکه هر دو تابع هموار هستند می توان از مدل های رگرسیونی شبه درستمایی نیمه پارامتری  $SPQR$  چیبو و مولر [۴، ۵]

<sup>5</sup>Score Equation

<sup>6</sup>Semi Parametric Quasi-Likelihood Regression Model

می شود. برای تقریب زدن مدل کامل با تابع پیوند بریده شده، تحت شرط

$$|g'(\cdot)|^2 \leq c \text{ مولر و استدمولر [۱۱] نشان دادند که}$$

$$\left\{ \int \left[ g(U_p + V_p) - g(U_p + s) \right] dF_{V_p|U_p}(s) \right\}^2$$

$$\leq \int g'(\xi)^2 (V_p - s)^2 dF_{V_p|U_p}(s)$$

$$\leq 2c \int (V_p^2 + s^2) dF_{V_p|U_p}(s)$$

و بنابراین:

$$E\left(\left(g(U_p + V_p) - g_p(U_p)\right)^2\right)$$

$$= E\left(\int \left[ g(U_p + V_p) - g(U_p + s) \right] dF_{V_p|U_p}(s) \right)^2$$

$$\leq 2c E\left(V_p^2 + E(V_p^2 | U_p)\right) = 2c E(V_p^2)$$

$$\leq 2c \sum_{j=p+1}^{\infty} \beta_j^2 \sum_{j=p+1}^{\infty} \sigma_j^2. \quad (2)$$

خطای تقریبی مدل بریده شده به طور مستقیم به  $\text{Var}(V_p)$  مربوط می شود و توسط فرمول  $\sigma_j^2 = \text{Var}(\varepsilon_j)$ ،  $j = 1, 2, \dots$  کنترل می شود. اکنون با در نظر گرفتن

$$\varepsilon_j^{(i)} = \int X_i(t) p_j(t) d\omega(t), \quad i = 1, \dots, n$$

مدل کامل با خطاهای استاندارد شده  $e'_i$  به صورت

$$Y_i = g\left(\alpha + \sum_{j=1}^{\infty} \beta_j \varepsilon_j^{(i)}\right) + e'_i \sigma\left(\alpha + \sum_{j=1}^{\infty} \beta_j \varepsilon_j^{(i)}\right),$$

است. با پیشگوهای خطی بریده شده  $\eta$  و میانگین های  $\mu$ ،

$$\eta_i = \alpha + \sum_{j=1}^p \beta_j \varepsilon_j^{(i)}, \quad \mu_i = g(\eta_i), \quad i = 1, \dots, n$$

و  $p$ -مدل بریده شده به صورت

$$Y_i^{(p)} = g_p\left(\alpha + \sum_{j=1}^p \beta_j \varepsilon_j^{(i)}\right) + e'_i \bar{\sigma}_p\left(\alpha + \sum_{j=1}^p \beta_j \varepsilon_j^{(i)}\right), \quad (3)$$

تعریف می شود، که در آن  $\bar{\sigma}_p$  و  $g_p$  ثابت هستند. در مدل بریده شده (۳) برای در نظر گرفتن عرض از مبدأ مجموع از یک قرار داده و فرض می شود  $\beta_0 = \alpha$  و  $\varepsilon_0^{(i)} = 1$

### ۳ برآورد پارامترهای مدل خطی

#### تابعی تعمیم یافته

برای برآورد پارامترها  $\beta' = (\beta_0, \dots, \beta_p)$  با ثابت نگه داشتن  $p$  از روش شبه درستمایی و دربرن [۱۵] باید معادله امتیاز<sup>۵</sup> به صورت

$$U(\beta) = 0 \quad (4)$$

استفاده کرد. چون تابع پیوند دلخواه است برای محدودیت‌های همواری و یکنوایی باید پارامترها و برآوردها

$$\|\beta\| = 1, \|\hat{\beta}\| = 1$$

باشند. قرار دهید  $\hat{\eta}_i = \sum_{j=0}^p \hat{\beta}_j \varepsilon_j^{(i)}$  و برآورد تابع پیوند به صورت  $\hat{g}(\cdot)$  و مشتق اول آن  $\hat{g}'(\cdot)$  از هموارسازی نمودار پراکنش  $(\hat{\eta}_i, Y_i)$  و برآورد تابع واریانس  $\hat{\sigma}^2(\cdot)$  با هموارسازی نمودار پراکندگی

$$(\hat{\mu}_i, \hat{\varepsilon}_i^2)_{i=1, \dots, n}$$

به دست می‌آید، که در آن  $\hat{\mu}_i = \hat{g}(\hat{\eta}_i)$  برآورد میانگین پاسخ و مجذور مانده‌ها  $\hat{\varepsilon}_i^2 = (Y_i - \hat{\mu}_i)^2$  است. اکنون با استفاده از معادله امتیاز نیمه پارامتری

$$U(\beta) = \sum_{i=1}^n (Y_i - \hat{g}(\eta_i)) g'(\eta_i) \varepsilon^{(i)} / \sigma^2(\hat{g}(\eta_i)) \quad (8)$$

برآوردها به صورت  $\hat{\beta} = (\hat{\beta}_0, \dots, \hat{\beta}_p)$  نتیجه می‌شود و در نهایت با حل معادله امتیاز (۵) و (۸) برآوردهای پارامترهای رگرسیونی تابعی به صورت

$$\hat{\beta}(t) = \hat{\beta}_0 + \sum_{j=1}^p \hat{\beta}_j \rho_j(t) \quad (9)$$

حاصل می‌شود. ماتریس‌های  $D$  و  $\Gamma$  برای حالت  $SPQR$  اصلاح و جایگزین برآوردها می‌شوند.

### ۱.۳ بازه اطمینان و آزمون فرضیه برای

#### پارامترهای رگرسیونی تابعی

یک تابع هسته انتگرال پذیر  $L^2$ ، به صورت  $\mathbb{R} \rightarrow \tau^2 : R(s, t)$  در نظر بگیرید و تعریف کنید  $A_R$  را به صورت یک عملگر انتگرال خطی

$$A_R : L^2(d\omega) \rightarrow L^2(d\omega)$$

روی فضای هیلبرت  $L^2(d\omega)$  برای  $f \in L^2(d\omega)$

$$(A_R f)(t) = \int f(s) R(s, t) d\omega(s). \quad (10)$$

عملگرهای  $A_R$  عملگرهای هیلبرت-اشمیت خودالحاقی فشرده<sup>۷</sup> هستند. اگر

$$\int |R(s, t)|^2 d\omega(s) d\omega(t) < \infty,$$

آنگاه مطابق با کانوی [۶] می‌تواند به صورت قطری شده باشند.

حالت‌های خاصی از عملگرهای انتگرال که در این مقاله مورد نظر هستند عملگر اتو کوواریانس  $A_K$  از  $X$  با هسته

$$k(s, t) = cov(X(s), X(t)) = E(X(s)X(t)), \quad (11)$$

و عملگر اتو کوواریانس تعمیم یافته  $A_G$  با هسته

$$G(s, t) = E \left( \frac{g'(\eta)^2}{\sigma^2(\mu)} X(s)X(t) \right), \quad (12)$$

عملگرهای هیلبرت-اشمیت هستند.  $A_R$  یک فضای متریک در  $L^2$  به صورت

$$\begin{aligned} d_R^2(f, g) &= \int (f(t) - g(t)) (A_R(f - g))(t) d\omega(t) \\ &= \int \int (f(s) - g(s)) (f(t) - g(t)) R(s, t) d\omega(s) d\omega(t) \end{aligned}$$

تولید می‌کند، که برای  $f, g \in L^2(d\omega)$  پایه یکامتامد دلخواه  $\{\rho_j, j = 1, 2, \dots\}$  هسته هیلبرت-اشمیت  $R$  را می‌توان به صورت

$$R(s, t) = \sum r_{kl} \rho_k(s) \rho_l(t),$$

برای ضرایب مناسب  $\{r_{kl}, k, l = 1, 2, \dots\}$  دانفورد و استوارتز [۷] بیان کرد.

با به کار بردن هر تابع مشخص  $h \in L^2$

$$h_{p,j} = \int h(s) \rho_j(s) d\omega(s),$$

و با نشان دادن توابع ویژه و مقادیر ویژه  $A_R$  با

$$\{\rho_j^R, \lambda_j^R, j = 1, 2, \dots\}$$

فاصله  $d_R$  را می‌توان به صورت

$$\begin{aligned} d_R^2(f, g) &= \sum_{k,l} r_{k,l} (f_{\rho,k} - g_{\rho,k}) (f_{\rho,l} - g_{\rho,l}) \\ &= \sum_k \lambda_k^R (f_{\rho,k} - g_{\rho,k})^2, \end{aligned} \quad (13)$$

بیان کرد. برای مدل  $p_n$ -بریده شده مربوط به مدل (۱) متر  $d_G$  به کار می‌رود که نسبت به متر  $L^2$  تحت شرایط ساده‌تری حدهای مجانبی محاسبه می‌شوند. این متر برای مدل بریده شده به صورت (۱)

$$\begin{aligned} d_G^2(\hat{\beta}, \beta) &= \int \int (\hat{\beta}(s) - \beta(s)) (\hat{\beta}(t) - \beta(t)) \\ &E \left( \frac{g'(\eta)^2}{\sigma^2(\mu)} X(s)X(t) \right) d\omega(s) d\omega(t), \end{aligned}$$

در نظر گرفته می‌شود که برای هر  $p$  با

$$d_{G,p}^2(\hat{\beta}, \beta) = (\hat{\beta} - \beta)' \Gamma (\hat{\beta} - \beta),$$

تقریب زده می‌شود. دقت کنید که اگر روش نیمه پارامتری  $SPQR$  به کار می‌رود، چون در این مدل توابع واریانس و پیوند نامعلوم هستند و به‌طور ناپارامتری برآورد می‌شوند، باید شرایطی در نظر گرفته شود که مهم‌ترین آن‌ها عبارت است از این که باید تابع پیوند  $g$  یکنواخت و معکوس پذیر مشتق پذیر باشد و برای هر  $c > 0$ ،  $\|g'(\cdot)\| \leq c$ ،  $\|g''(\cdot)\| \leq c$ . تابع واریانس  $\sigma^2(\cdot)$  به‌طور پیوسته مشتق پذیر است و  $\delta > 0$  وجود دارد، به طوری که  $\delta > \sigma(\cdot)$ . همچنین در دنباله‌ای از مدل‌های تقریبی بریده شده  $p_n$ ، با تعداد

<sup>7</sup>Compact Self-adjoint Hilbert-schmidt

در حالی که توابع پیوند و واریانس نامعلوم هستند در رهیافت  $SPQR$ ، برای به دست آوردن برآورد ناپارامتری از توابع و مشتق‌ها، هموارسازی نمودارهای پراکنش اعمال می‌شود و سپس برآورد پارامتری  $\hat{\beta}$  از حل معادله امتیاز نیمه پارامتری (۸) به دست می‌آید. بعد از چند تکرار برآورد ناپارامتری تابع پیوند  $\hat{g}$  و مشتق اولش  $\hat{g}'$  و تابع واریانس  $\hat{\sigma}^2$  به دست می‌آیند. این برآوردگرهای منحنی ناپارامتری توابع با هموارکننده‌های هسته درجه دوم یا خطی موضعی با یک پهنای باند  $h$  در هموارسازی به دست می‌آیند. مولر و استدمولر، [۱۱] نشان دادند تحت شرایطی توزیع تقریبی (۱۴) برای رهیافت  $SPQR$  نیز برقرار است.

یک مسئله رایج در مدل‌های رگرسیونی به این صورت است که

$$H_0: \beta = \beta_0 \text{ ثابت}$$

که حالت خاصی از آزمون  $H_0: \beta \equiv \beta_0$  برای یک تابع پارامتری مشخص  $\beta_0$  است. با در نظر گرفتن  $\beta_0(t) = \sum \beta_{0j} \rho_j(t)$ ، فرض  $H_0: \beta_j = \beta_{0j}$ ،  $j = 0, 1, 2, \dots$  می‌شود و  $H_0$  رد می‌شود وقتی آماره‌ی آزمون بیان‌شده در (۱۴) بیشتر از مقدار بحرانی  $\Phi(1 - \alpha)$  باشد.

در عمل بازه‌های اطمینان برای تابع پارامتر رگرسیونی  $\beta$  نیز موردنظر است. در حالی که  $p = p_n$  معلوم است و برآوردهای  $\hat{\beta}$  برای بردار  $\beta$  مشخص شده‌اند یک ناحیه اطمینان مجانبی  $(1 - \alpha)$  مطابق با (۱۴) با  $c(\alpha)$  مشخص می‌شود که در آن

$$c(\alpha) = [p + 1 + \sqrt{2(p+1)\Phi(1 - \alpha)}] / n,$$

$\Gamma$  معمولاً با  $\bar{\Gamma}$  یا  $\hat{\Gamma}$  جایگزین می‌شود. اکنون فرض کنید  $(e_1, \lambda_1), \dots, (e_{p+1}, \lambda_{p+1})$  بردارهای ویژه و مقادیر ویژه ماتریس  $\Gamma$  باشند، و قرار دهید

$$\omega_k(t) = \sum_{l=1}^{p+1} \rho_l(t) e_{kl}, \quad k = 1, \dots, p+1, \quad e_k = (e_{k1}, \dots, e_{k,p+1})'$$

آنگاه برای  $n$  بزرگ و  $p$  یک بازه اطمینان تقریبی هم‌زمان  $(1 - \alpha)$  برای تابع  $\beta(t)$  به صورت

$$\hat{\beta}(t) = \pm \sqrt{c(\alpha) \sum_{k=1}^{p+1} \frac{\omega_k(t)^2}{\lambda_k}} \quad (17)$$

است. در عمل بازه اطمینان هم‌زمان با جایگزین کردن برآوردهای  $\omega_k$  و  $\lambda_k$  که از ماتریس  $\bar{\Gamma}$  یا  $\hat{\Gamma}$  به جای  $\Gamma$  به دست آمده‌اند، محاسبه می‌شود.

جملات پیشگویی  $p_n$ ، اگر  $p_n \rightarrow \infty$  و  $n \rightarrow \infty$  آنگاه  $p_n n^{-1/4} \rightarrow 0$  (برای جزئیات شرایط مدل  $SPQR$  به مولر و استدمولر، [۱۱] مراجعه شود). اکنون با در نظر گرفتن فرض‌های اصلی و شرط‌های مدل  $SPQR$ ، می‌توان نشان داد، اگر  $n \rightarrow \infty$

$$\frac{n(\hat{\beta} - \beta)' \Gamma_{p_n} (\hat{\beta} - \beta) - (p_n + 1)}{\sqrt{2(p_n + 1)}} \xrightarrow{d} N(0, 1), \quad (14)$$

که مولر و استدمولر، [۱۱] با در نظر گرفتن  $p_n = p$ ،  $\bar{\Gamma} = \frac{1}{n} E(D'D)$  را جایگزین نمودند. در صورتی که فقط پارامترهای شیب  $\beta_1, \beta_2, \dots$  (بدون عرض از مبدأ) موردنظر هستند  $p_n$  با  $p_n - 1$  و به جای ماتریس  $(p+1) \times (p+1)$ ،  $\Gamma$  از ماتریس  $p \times p$  استفاده می‌شود که از حذف اولین سطر و ستون به دست آمده است. برای همگرایی تابع پارامتر برآورد شده  $\hat{\beta}(\cdot)$ ، فاصله  $d_G$  و بسط به صورت

$$\hat{\beta}(t) = \sum \hat{\beta}_{\rho_j^G} \rho_j^G(t)$$

مربوط به تابع پارامتر  $\hat{\beta}(\cdot)$  در پایه  $\{\rho_j^G, j = 1, 2, \dots\}$  باید پایه ویژه عملگر  $AG$  متناسب با مقادیر ویژه  $\lambda_j^G$  به کار گرفته شود. بنابراین

$$\begin{aligned} d_G^2(\hat{\beta}(\cdot), \beta(\cdot)) &= \int \int (\hat{\beta}(s) - \beta(s)) G(s, t) \\ &\quad (\hat{\beta}(t) - \beta(t)) d\omega(s) d\omega(t) \\ &= \sum_{j=1}^p \lambda_j^G (\hat{\beta}_{\rho_j^G} - \beta_{\rho_j^G})^2 + \sum_{j=p+1}^{\infty} \lambda_j^G \beta_{\rho_j^G}^2 \\ &= (\hat{\beta}^G - \beta^G)' \Gamma^G (\hat{\beta}^G - \beta^G) + \sum_{j=p+1}^{\infty} \lambda_j^G \beta_{\rho_j^G}^2. \end{aligned}$$

که در آن  $\beta^G = (\beta_{\rho_1^G}, \dots, \beta_{\rho_p^G})'$ ،  $\hat{\beta}^G = (\hat{\beta}_{\rho_1^G}, \dots, \hat{\beta}_{\rho_p^G})'$  و  $\Gamma^G$  ماتریس قطری است.  $\varepsilon_j^G$  به صورت

$$\varepsilon_j^G = \frac{g'(\eta)}{\sigma(\mu)} \int X(t) \rho_j^G(t) d\omega(t),$$

مشخص می‌شود، به طوری که

$$\begin{aligned} E(\varepsilon_j^G \varepsilon_k^G) &= \int \int G(s, t) \rho_j^G(s) \rho_k^G(t) d\omega(s) d\omega(t) \\ &= \delta_{ij} \lambda_j^G. \end{aligned} \quad (15)$$

**نتیجه ۱.۳.** اگر تابع پارامتری  $\beta(\cdot)$  ویژگی

$$\sum_{j=p+1}^{\infty} E(\varepsilon_j^{G(n)}) \left[ \int \beta(t) \rho_j^G(t) d\omega(t) \right]^2 = o\left(\frac{\sqrt{p_n}}{n}\right), \quad (16)$$

را داشته باشد، آنگاه وقتی  $n \rightarrow \infty$

$$\frac{n \int \int (\hat{\beta}(s) - \beta(s)) (\hat{\beta}(t) - \beta(t)) G(s, t) d\omega(s) d\omega(t) - (p_n + 1)}{\sqrt{2p_n + 1}} \xrightarrow{d} N(0, 1)$$

(برای اثبات به مولر و استدمولر، [۱۱] مراجعه شود)

### ۲.۳ برآورد درستنمایی تاوانیده

برای برآورد درستنمایی تاوانیده معمولاً تابع تاوان<sup>۸</sup> که توسط مایر و همکاران [۱۰] ارائه شده است به صورت

$$P_{\lambda, \varphi}(\beta) = \lambda (\|\beta\|^2 + \varphi \|\beta''\|^2)^{\frac{1}{2}}$$

در نظر گرفته می شود، که در آن  $\|\beta\|^2 = \int_{D_1} \{\beta(t)\}^2 dt$  نرم  $L^2$  از  $\beta$  و  $\beta''(t) = \frac{\partial^2 \beta(t)}{\partial t^2}$  است.  $\lambda, \varphi$  پارامترهای تاوان غیرمنفی هستند که به ترتیب تنگی و همواری را کنترل می کنند. برای به دست آوردن برآورد ماکسیمم درستنمایی تاوانیده تابع  $\beta(\cdot)$  باید عبارت

$$\frac{1}{n} \sum_{i=1}^n \ell(y_i, \beta(t) X_i(t)) - \sum_{j=1}^p P_{\lambda, \varphi}(\beta_j(t))$$

مینیمم شود به طوری که با در نظر گرفتن

$$\beta(t) = \sum_{j=1}^p \beta_j \rho_j(t)$$

تاوان به صورت

$$P_{\lambda, \varphi}(\beta) = \lambda (\rho'(\Psi + \varphi \Omega) \rho)^{\frac{1}{2}}$$

می شود، که در آن  $\Psi$  ماتریس  $p \times p$  که  $(j, k)$  امین عنصر آن

$$\Psi_{j,k} = \int \beta_j(t) \beta_k(t) dt, \quad j, k = 1, \dots, p$$

و  $\Omega$  یک ماتریس  $p \times p$  که  $(j, k)$  امین عنصر آن

$$\Omega_{j,k} = \int \beta_j''(t) \beta_k''(t) dt, \quad j, k = 1, \dots, p$$

هستند. در صورتی که تاوان به صورت

$$P_{\lambda, \varphi}(\beta) = \lambda (\rho' k \rho)^{\frac{1}{2}}$$

در نظر گرفته شود نوعی تاوان لاسو<sup>۹</sup> است که در آن  $\Psi + \varphi \Omega = k$  یک ماتریس  $p \times p$  متقارن معین مثبت است. در صورتی که  $k \varphi = I$  تاوان لاسو معمولی است.

### ۴ انتخاب رتبه مدل

همانطور که در قبل اشاره شد یکی از مطالب مهم در مدل بریده شده خطی تابعی تعمیم یافته انتخاب مقدار مناسب برای  $p$  یا تعداد توابع ویژه ای که برای برازش تابع  $\beta(t)$  به کار می روند، به عبارت دیگر انتخاب رتبه برای مدل

است. یک رهیافت برای این انتخاب استفاده از معیار آکائیکه<sup>۱۰</sup>  $AIC$  است که تمایل به افزایش رتبه مدل دارد. فرض کنید بردار پیشگوی خطی  $\eta_p$  شامل  $n$  مؤلفه  $\eta_{p,i} = \sum_{j=1}^p \varepsilon_j^i \beta_j$ ،  $i = 1, \dots, n$  باشد و بردار  $\hat{\eta}_p$  شامل مؤلفه های  $\hat{\eta}_{p,i} = \sum_{j=1}^p \varepsilon_j^i \hat{\beta}_j$  شامل مؤلفه های  $G$  را مشتق تابع معکوس پیوند  $g$  در نظر بگیرید بنابراین  $Y$  دارای چگالی

$$f_Y(y) = \exp(y\eta + a(y) - G(\eta))$$

است. در عمل  $\tilde{\sigma}^2(\eta) = g'(\eta)$  می باشد. انحراف<sup>۱۱</sup> به صورت

$$\mathcal{D} = -2\ell_n(Y, \hat{\eta}_p) + 2\ell_n(Y, g^{-1}(Y)),$$

تعریف می شود، که در آن لگاریتم شبه درستنمایی

$$\ell_n(Y, \hat{\eta}_p) = \sum_{i=1}^n Y_i \hat{\eta}_{i,p} - \sum_{i=1}^n G(\hat{\eta}_{i,p}).$$

است. با استفاده از بسط تیلور

$$\begin{aligned} -2\ell_n(Y, \hat{\eta}_p) &= -2\ell_n(Y, \eta_p) \\ &+ 2 \left( \nabla_{\beta_p} \ell_n(Y, \hat{\eta}_p) \right)' (\beta_p - \hat{\beta}_p) \\ &+ (\beta_p - \hat{\beta}_p)' \left( \frac{\partial^2}{\partial \beta_k \partial \beta_l} \ell_n(Y, \hat{\eta}_p) \right) (\beta_p - \hat{\beta}_p) \end{aligned}$$

که در آن با توجه به معادله امتیاز، عبارت دوم سمت راست صفر و ماتریس به صورت فرم درجه دوم  $(D'D)$  است. طبق مولر و استدمولر [۱۱]. مقدار امید ریاضی فرم درجه دوم

$$n(\beta_p - \hat{\beta}_p)' \left( \frac{D'D}{n} \right) (\beta_p - \hat{\beta}_p),$$

به طور تقریبی مقدار  $p$  است. چون

$$\begin{aligned} -2\ell_n(Y, \eta_p) &= -2\ell_n(Y, \eta) - 2 \sum_{i=1}^n (Y_i - g(\eta_i)) (\eta_{i,p} - \eta_i) \\ &+ \sum_{i=1}^n g'(\eta_i) (\eta_{i,p} - \eta_i)^2 \end{aligned} \quad (18)$$

در نتیجه

$$\begin{aligned} E(\mathcal{D}) &= n \sum_{k,l=p+1} E \left( g'(\eta) \varepsilon_k \varepsilon_l \right) \beta_k \beta_l - p(1 + o(1)) + E_n \\ &= n \sum_{k,l=p+1} E \left( \frac{g''(\eta)}{\sigma^2(\eta)} \varepsilon_k \varepsilon_l \right) \beta_k \beta_l - p(1 + o(1)) + E_n, \end{aligned}$$

که در آن عبارت  $E_n$  به  $p$  وابسته نیست. مولر و استدمولر [۱۱] با به کار گیری قانون اعداد بزرگ نشان دادند که اگر  $p$  در بازه  $(p_0, cn^{\frac{1}{2}})$  قرار بگیرد آنگاه

$$\mathcal{D}/E(\mathcal{D}) \xrightarrow{P} 1.$$

<sup>8</sup>Penalty Function

<sup>9</sup>Lasso

<sup>10</sup>Akaike Information Criterion

<sup>11</sup>Deviance

همچنین آن‌ها نشان دادند

$$\begin{aligned} d(\hat{\beta}(\cdot), \beta(\cdot)) &= \int \int (\hat{\beta}(s) - \beta(s)) G(s, t) (\hat{\beta}(t) - \beta(t)) d\omega(s) d\omega(t) \\ &= (\hat{\beta}_p - \beta_p)' \Gamma(\hat{\beta}_p - \beta_p) + \sum_{k,j=p+1}^{\infty} \gamma_{j,k} \beta_j \beta_k \\ &+ 2 \sum_{j=1}^p \sum_{k=p+1}^{\infty} \gamma_{j,k} (\hat{\beta}_j - \beta_j) \beta_k \end{aligned}$$

که در آن

$$\gamma_{k,l} = E \left( \frac{g^{\vee}(\eta)}{\sigma^{\vee}(\eta)} \varepsilon_k \varepsilon_l \right).$$

اکنون

$$E(d(\hat{\beta}(\cdot), \beta(\cdot))) = p/n(1 + o(1)) + \sum_{k,j=p+1}^{\infty} \gamma_{j,k} \beta_j \beta_k (1 + o(1)).$$

این استنباط نشان می‌دهد که تابع هدف  $(d(\hat{\beta}(\cdot), \beta(\cdot)))$  با کمینه ساختن به‌طور مجانبی به  $E(\mathcal{D}/n) + 2p/n$  نزدیک می‌شود. در صورتی که دقت تابع هدف و  $AIC$  مطرح باشد، اثبات می‌شود که رتبه  $p_A$  با  $AIC$  و  $p_d$  با مینیم کردن تابع هدف انتخاب می‌شوند. آنگاه اگر  $an \rightarrow \infty$  و  $p_d/p_A \rightarrow 1$  در عمل  $AIC$  و ملاک اطلاع بیزی  $^{12} BIC$  را با به دست آوردن اولین انحراف یا شبه‌انحراف  $\mathcal{D}(p)$ ، وابسته به مرتبه مدل  $p$  بکار می‌بریم. که در شبه‌درست‌نمایی یا ماکسیمم درست‌نمایی اگر با تابع پیوند معلوم باشد به راحتی امکان‌پذیر است و در صورتی که نامعلوم باشد، لگاریتم درست‌نمایی و رهیافت  $SPQR$  با تابع پیوند نامعلوم نیازمند انتگرال‌گیری است. در نهایت عبارت

$$C(p) = \mathcal{D}(p) + \mathcal{P}(p) \quad (19)$$

کمینه می‌شود که در آن  $\mathcal{P}$  جریمه است، به طوری که  $2p$  برای  $AIC$  و  $p \log n$  برای  $BIC$  انتخاب می‌شود. روش‌های دیگری برای تعیین  $p$  وجود دارد که می‌توان به استفاده از اعتبارسنجی متقابل رایس و سیلورمن، [۱۴] و اختلاف نسبی بین معیار پیرسون و انحراف چپو و مولر، [۴] اشاره نمود.

## ۵ بررسی چند مثال

در این بخش یک مطالعه شبیه‌سازی و دو مثال کاربردی موردبررسی قرار گرفته است.

### ۱.۵ شبیه‌سازی

در یک مدل خطی تابعی تعمیم‌یافته علاوه بر انتخاب  $p$ ، تعیین یک پایه یک‌معامد  $\{\rho_j, 1, 2, \dots\}$  مهم است. معمولاً دو انتخاب داریم، یا از یک پایه استاندارد ثابت مانند پایه فوریه  $\rho_j \equiv \sqrt{2} \sin(\pi j t)$ ،  $\varphi_j \equiv \varphi_j \equiv \sqrt{2} \sin(\pi j t)$ ،  $t \in [0, 1]$ ،  $j \geq 1$  استفاده می‌شود، یا توابع ویژه عملگر کوواریانس  $A_K$  که در بخش ۴ در روابط (۱۱) و (۱۰) دیدیم. در اینجا از الگوریتم دوم استفاده می‌شود و الگوریتمی که

برای برآورد توابع ویژه در کاپرا و مولر [۱]؛ رایس و سیلورمن [۱۴] ارائه شده، به کار می‌رود.  $p$  مؤلفه مدل تعیین می‌شود و بنابراین  $i$  امین واحد مشاهده شده فرآیند به  $p$  پیشگویی کننده کاهش می‌یابد

$$\varepsilon_j^{(i)} = \int X_i(t) \rho_j(t) d\omega(t), \quad j = 1, \dots, p.$$

توابع ویژه برآورد شده به جای  $\rho_j$  جایگزین و انتگرال‌ها به‌طور عددی محاسبه می‌شوند. پارامترهای  $\alpha$  و  $\beta_1, \dots, \beta_p$  در مدل تابعی تعمیم‌یافته با حل معادله امتیاز برآورد می‌شوند. الگوریتم کمترین مربعات بحث شده در مک کلا و نلدر [۹] برای حالتی از یک مدل خطی تعمیم‌یافته یا شبه درست‌نمایی با تابع پیوند معلوم و الگوریتم  $^{13} QLUE$  مطرح شده در چپو و مولر [۵] برای مدل  $SPQR$  با تابع پیوند نامعلوم تکرار می‌کنیم. هدف از ارائه این مثال شبیه‌سازی مونت کارلویی مقایسه  $AIC$  و  $BIC$  برای تعیین رتبه  $p$  است و بررسی توان آزمون‌های آماری اثر رگرسیون در مدل رگرسیون تابعی تعمیم‌یافته و همچنین بررسی رفتار فرآیند نیمه پارامتری  $SPQR$  برای رگرسیون تابعی، در مقایسه با پیاده‌سازی شبه درست‌نمایی و ماکزیمم با تابع پیوند مشخص است. شبیه‌سازی به این صورت است که فرآیندهای شبه تصادفی مبنی بر  $2^0$  تابع اول پایه فوریه  $X(t) = \sum_{j=1}^{2^0} \varepsilon_j \varphi_j(t)$  با استفاده از متغیرهای شبه تصادفی نرمال  $\varepsilon_j \sim N(0, 1/j^2)$ ،  $j \geq 1$  تولید شده‌اند. فرض کنید  $\beta_j = 1/j$ ،  $1 \leq j \leq 3$ ،  $\beta_0 = 0$ ،  $\beta_j = 0$  و  $j > 3$ . فرض کنید

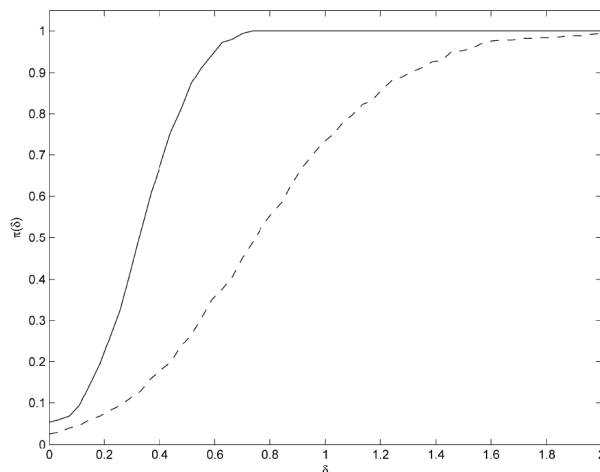
$$\beta(t) = \sum_{j=1}^{2^0} \beta_j \varphi_j(t)$$

$$P(X(\cdot)) = g(\beta_0 + \sum_{j=1}^{2^0} \beta_j \varepsilon_j)$$

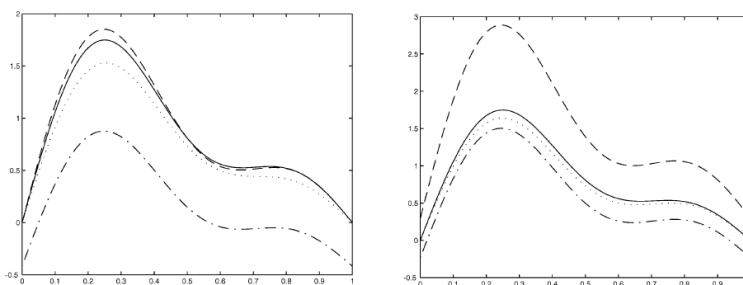
دو پیوند، پیوند لوجیت به طوری که  $[g(x) = \exp(x)/(1 + \exp(x))]$  و پیوند لگ-لگ-مکمل  $(c - \log \log)$ ،  $[g(x) = \exp(-\exp(-x))]$  در نظر گرفته می‌شود. آنگاه  $Y|X \sim \text{Binomial}(P(X), 1)$  به صورت یک متغیر تصادفی شبه‌برنولی با احتمال  $P(X)$  و نمونه  $(X_i(t), Y_i)$ ،  $i = 1, \dots, n$ ، تولید می‌شود. برای نشان دادن نتایج تقریبی که در (۱۴) دیدیم، توابع توان تجربی برای داده‌های تولیدشده و با به کار بردن تابع پیوند لوجیت و آماره آزمون  $T$  که در سمت چپ رابطه (۱۵) برای آزمون کردن فرضیه  $H_0: \beta_j = 0$ ،  $j = 1, 2, \dots$  یعنی عدم وجود اثر رگرسیونی دیدیم، محاسبه می‌شوند. این آزمون یک آزمون یک‌طرفه در سطح ۵٪ است، به طوری که اگر  $|T| > \Phi^{-1}(0.95)$  فرض  $H_0$  رد می‌شود. با تولید ۵۰۰ نمونه مونت کارلو نرخ رد برای حجم نمونه  $n = 50, 200$  به صورت تابع  $\delta$ ،  $2 \leq \delta \leq 0$  تعیین شد که بردار پارامتر شرح داده‌شده در پاراگراف قبلی، در  $\delta$  ضرب شده است و به وسیله  $(\delta/3, \delta/2, \delta)$  تعیین می‌شود. تابع توان برای دو حجم نمونه ۵۰، ۲۰۰ در شکل ۱ رسم شده است. این شکل نشان می‌دهد که اندازه نمونه نقش مهمی توان آزمون دارد.

<sup>12</sup>Bayesian Information Criterion

<sup>13</sup>Quasi Likelihood with Unknown linke and variance function Estimation



شکل ۱: تابع توان آزمون معناداری رگرسیون لوژیستیک تابعی در سطح ۵٪ بر اساس ۵۰۰ شبیه‌سازی، برای نمونه‌ای به اندازه ۵۰ (خط چین) و ۲۰۰ (خط ممتد)، با  $p = 3$



شکل ۲: شکل سمت راست داده‌ها از پیوند لگ-لگ مکمل و شکل چپ از تابع پیوند لوجیت تولید شده است. خط ممتد در هر دو نمودار خط هدف است و برآوردها با فرض پیوند لوجیت با (خط چین)، پیوند با (خط چین-نقطه) و روش لگ-لگ مکمل  $SPQR$  با (نقطه چین) مشخص شده است.

۳. رهیافت نیمه پارامتری  $SPQR$  با پیوند نامعلوم

محاسبه شد. با توجه به شکل ۲ درمی‌یابیم که روش  $SPQR$  یک رهیافت مفید است. بنابراین نتیجه می‌شود که انتخاب تابع پیوند یک مسئله مهم در این مدل‌ها می‌باشد و انتخاب نادرست برای آن می‌تواند باعث بروز مشکلاتی در استنباط‌ها شود.

برای این که خوب بودن روش نیمه پارامتری  $SPQR$  بررسی شود، میانگین

برآورد تابع پارامتر  $(\cdot)\hat{\beta}$  با ۵۰ بار اجرای مونت کارلو در حالات زیر

۱. استفاده از تابع پیوند لوجیت

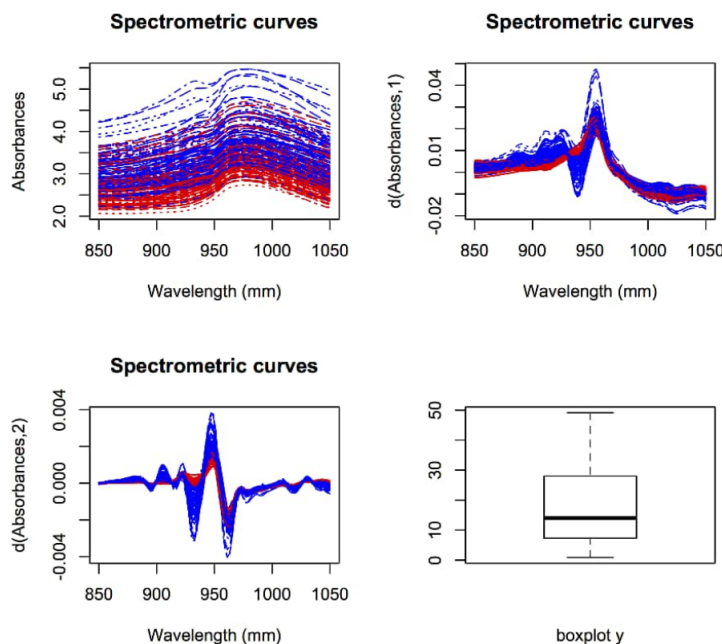
۲. استفاده از تابع پیوند لگ-لگ مکمل

## ۲.۵ مثال کاربردی ۱

شده در ۱۰۰ طول موج بین ۸۵۰ تا ۱۰۵۰ میلی‌متر به عنوان متغیر کمکی تابعی فرض می‌شود. این مجموعه داده بخشی از مجموعه داده اصلی است که در <http://stat.cmu.edu> در دسترس است. نمودار متغیر تابعی میزان جذب و مشتق اول و دوم و نمودار جعبه‌ای متغیر پاسخ قبل از رسته‌ای کردن این متغیر در شکل ۳ مشاهده می‌شود.

فرض کنید در یک بررسی متغیر پاسخ چربی گوشت است که برای چربی کمتر از ۱۵ عدد صفر و برای چربی بیشتر از ۱۵ عدد یک در نظر گرفته شده است و متغیر پاسخ به صورت صفر و یک است. میزان جذب اندازه‌گیری





شکل ۳: رنگ آبی مربوط به چربی کمتر از ۱۵ و رنگ قرمز مربوط به چربی بیشتر از ۱۵ است.

مدل لگ-لگ مکمل تابع پیوند به صورت

$$P_i = P(Y = 1 | X_i(t), t \in \tau) = 1 - \exp(-\exp(\alpha + \int \beta(t)X_i(t)dt))$$

است. مقدار  $AIC$  سه مدل لوجیت، پروبیت و لگ-لگ مکمل به ترتیب ۲۷۰۳۲، ۲۷۱۹۴ و ۲۷۰۳۲ به دست آمد که مدل لگ-لگ مکمل با توجه به  $AIC$  کمتر مدل مناسب است. همچنین نرخ طبقه‌بندی نادرست برای مدل لوجیت، پروبیت، لگ-لگ مکمل و  $SPQR$  به ترتیب ۳۵٪، ۳۳٪ و ۳۳٪ و ۳۸٪ به دست آمد که نشان می‌دهد که نتایج مدل لگ-لگ مکمل به رهیافت  $SPQR$  نزدیک‌تر است و برای این داده‌ها از دقت بیشتری برخوردار است.

بلند هستند. با انتخاب پیوند لوجیت، یک رگرسیون تابعی لوجیت به داده‌ها برازش می‌شود. خط سیر تولیدمثل برای مگس‌های با طول عمر بلند و عمر کوتاه به صورت جداگانه در شکل ۴ ترسیم شده است. تفاوت بین دو گروه به دلیل ازدحام منحنی‌های رسم شده به صورت چشمی قابل تشخیص نیست. برای انتخاب تعداد مؤلفه‌های مدل ملاک  $AIC$  به کار می‌رود. همان‌طور که در شکل ۵ مشاهده می‌شود، ملاک  $AIC$  به مرتبه مدل  $p$  وابسته است و مقدار مینیمم در ۶ است. بنابراین  $p = 6$  انتخاب می‌شود. برای انتخاب مقدار مناسب  $p = 6$  از معیار خطای پیشگویی اعتبارسنجی متقابل  $PE = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{p}_i^{(-i)})^2$  با آوردن  $p_i$  با خروج  $\hat{p}_i^{(-i)}$  از

برای تحلیل داده‌ها با استفاده از مدل خطی تابعی تعمیم‌یافته از بسته نرم‌افزاری `fda.usc` و تابع `fregre.glm` در نرم‌افزار `R` استفاده شد. با توجه به اینکه متغیر پاسخ دودویی است می‌توان از توابع پیوند لوجیت<sup>۱۴</sup> و پروبیت<sup>۱۵</sup> و لگ-لگ مکمل (`cloglog`) استفاده کرد که می‌دانیم برای پیوند لوجیت

$$P_i = P(Y = 1 | X_i(t), t \in \tau) = \frac{\exp(\alpha + \int \beta(t)X_i(t)dt)}{1 + \exp(\alpha + \int \beta(t)X_i(t)dt)}$$

برای پیوند پروبیت

$$P_i = P(Y = 1 | X_i(t), t \in \tau) = \Phi(\alpha + \int \beta(t)X_i(t)dt)$$

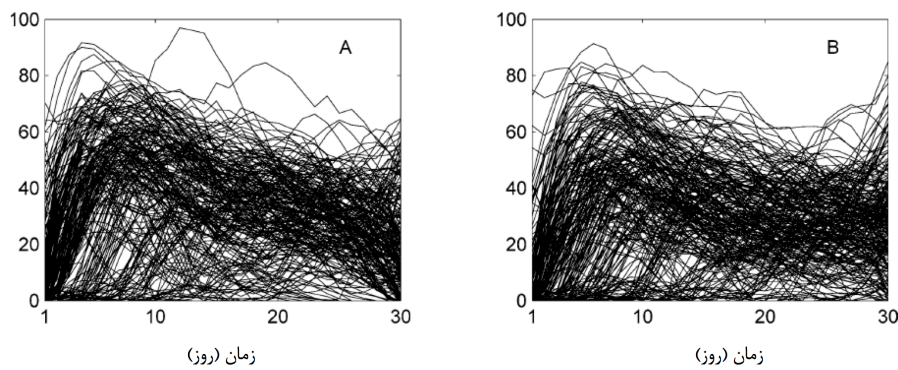
که در آن  $\Phi$  تابع توزیع تجمعی نرمال استاندارد است استفاده می‌شود. برای

### ۳.۵ مثال کاربردی ۲

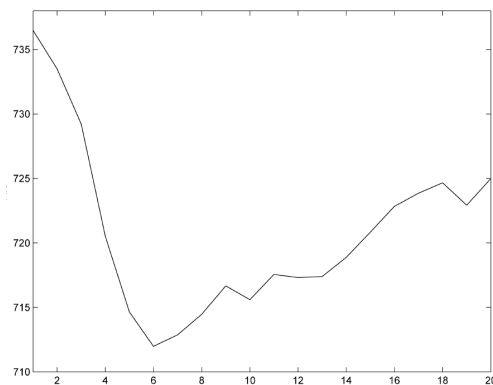
از هزار مگس مدیترانه‌ای مطرح شده در کری و همکاران [۳]، مگس‌هایی که ۳۴ روز بعد زنده مانده‌اند را انتخاب می‌کنیم و نمونه‌ای از ۵۳۴ مگس مدیترانه‌ای فراهم و برای پیش‌بینی، خط مسیر تخم‌گذاری از ۰ تا ۳۰ روز برای فرآیندهای پیشگویی  $X_i(t), t \in [0, 30], i = 1, \dots, 534$  استفاده می‌شود. اگر عمر مگس در ۳۰ روز گذشته ۱۴ روز یا بیشتر باشد در طبقه‌بندی طول عمر بلندمدت قرار می‌گیرد، در غیر این صورت طول عمر کوتاه دارد. از  $n = 534$  مگس، ۲۵۶ مگس طول عمر کوتاه و ۲۷۸ مگس دارای طول عمر

<sup>14</sup>Logit

<sup>15</sup>Probit



شکل ۴: پیش‌بینی خط سیر، مطابق با منحنی‌های تخم‌گذاری روزانه، برای  $n = 534$  مگس. (چپ) تولیدمثل ۲۵۶ مگس با طول عمر کوتاه و (راست) ۲۷۸ مگس با طول عمر بلند را نمایش می‌دهد.



شکل ۵: معیار اطلاعات آکائیکه ( $AIC$ ) به عنوان تابعی از مقادیر مؤلفه‌های مدل  $p$  برای داده مگس مدیترانه‌ای

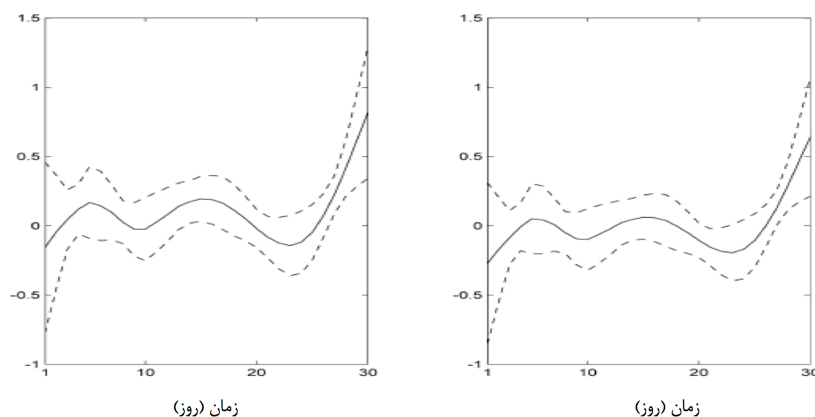
می‌کند که دوره‌های تناوب تخم‌گذاری با افزایش طول عمر مرتبط است. تحت تابع پیوند لوجیت، می‌دانیم احتمال طبقه‌بندی مگس با طول عمر بلند برای برآورد تابع پیوند ناپارامتری  $g(\eta) = \exp(\eta) / (1 + \exp(\eta))$  است. برای رهیافت  $SPQR$  استفاده شد، که هموارسازی خطی موضعی و پهنای باند ۰.۵۵ برای مراحل هموارسازی انتخاب شد. برای دو تابع پیوند، بزرگ‌ترین پیشگوی خطی  $\eta$ ، سپس بزرگ‌ترین مقادیر تابع پارامتر  $\beta(\cdot)$ ، با احتمال افزایش طول عمر مرتبط هستند. چون برآورد تابع پارامتر بین ۱۷ - ۱۲ روز و پس از ۲۶ روز نسبتاً بزرگ است، نتیجه می‌گیریم که در طول این دوره‌ها فعالیت باروری سنگین با افزایش طول عمر در ارتباط است. در مقابل، افزایش تولیدمثل بین ۱۲ - ۸ روز و ۲۶ - ۲۰ روز با کاهش طول عمر در ارتباط است.

مورد بررسی قرار گرفت. برای این مدل آزمون فرضیه و بازه اطمینان مجانبی برای پارامترهای رگرسیونی تابعی به دست آورده شد. مدل نیمه‌پارامتری  $SPQR$  برای حالتی که تابع پیوند نامعلوم است بررسی و در نهایت مدل و روش‌ها بر روی چند مثال پیاده‌سازی شد. همان‌طور که در مثال‌ها مشاهده شد

است، استفاده شد. از این معیار هم  $p = 6$  به دست آمد. برآورد نرخ طبقه‌بندی نادرست برای گروه مگس‌های طول عمر بلند، ۳۷٪ با پیوند لوجیت و ۳۵٪ برای پیوند  $SPQR$  است. درحالی‌که برای گروه مگس‌های کوتاه عمر ۴۷٪ برای لوجیت و ۴۸٪ برای  $SPQR$  است. بنابراین طبقه‌بندی صحیح مگس‌های کوتاه عمر را به سختی نشان می‌دهد. توابع پارامتر  $\hat{\beta}(\cdot)$  برای رگرسیون تابعی لوژستیک و  $SPQR$  با بازه‌های اطمینان هم‌زمان در شکل ۶ نشان داده شده است؛ مشاهده می‌شود که برآورد با پیوند ناپارامتری کاملاً نزدیک به برآورد با پیوند لوژستیک است، بنابراین در این حالت پیوند لوژستیک را انتخاب می‌کنیم. از بازه‌های اطمینان مجانبی نتیجه می‌گیریم تابع پیوند شیب صعودی را در سمت راست برای عمر ۳۰ روزه دارد، و شکل تابع پارامتر  $\hat{\beta}(\cdot)$  مشخص

## بحث و نتیجه‌گیری

یک مدل خطی تابعی تعمیم‌یافته بریده‌شده برای بررسی ارتباط بین یک متغیر پاسخ اسکالر گسسته و متغیرهای کمکی تابعی با رهیافت شبه درستنمایی



شکل ۶: برآورد پارامتر تابعی  $\hat{\beta}(\cdot)$  (توپر) و بازه‌های اطمینان هم‌زمان، (چپ) پیوند لوجیت، (راست) رهیافت ناپارامتری SPQR.

معمولاً رهیافت SPQR برای تحلیل مدل خطی تابعی تعمیم‌یافته مناسب است. رهیافت ناپارامتری رهیافت SPQR به کار گرفته شود و در نهایت بهترین مدل با توجه به داده‌ها می‌توان توابع پیوند مختلف در نظر گرفته شود و برای اطمینان انتخاب شود.

## مراجع

- [1] Capra, W. B. and Muller, H. G. (1997). An accelerated time model for response curves. *Journal of the American Statistical*, **92**, 72–83.
- [2] Cardot, H., Ferraty, F. and Sarda, P. (1999). Functional linear model. *Statistics & Probability Letters*, **45**, 11–22.
- [3] Carey, J. R., Liedo, P., Muller, H. G., Wang, J. L. and Chiou, J. M. (1998). Relationship of age patterns of fecundity to mortality, longevity and lifetime reproduction in a large cohort of Mediterranean fruit fly females. *Journal of the Gerontology: Biological Sciences*, **53A**, B245–B251.
- [4] Chiou, J. M. and Muller, H. G. (1998). Quasi-likelihood regression with unknown link and variance functions. *Journal of the American Statistical Association*, **93**, 1376–1387.
- [5] Chiou, J. M. and Muller, H. G. (1999). Nonparametric quasi-likelihood. *Journal of the American Statistical*, **27**, 36–64.
- [6] Conway, J. B. (1990). *A Course in Functional Analysis*, 2nd ed. Springer, New York.
- [7] Dunford, N. and Schwartz, J. T. (1963). *Linear operators. II. Spectral Theory*. Wiley, New York.
- [8] Faraway, J. J. (1997). Regression analysis for a functional response. *Technometrics*, **39**, 254–261.
- [9] McCullagh, P. and Nelder, J. A. (1989). *Generalized linear models*, 2nd ed. Chapman and Hall, London.
- [10] Meier, L, Van de Geer, S & Bühlmann, P (2009). High-dimensional additive modeling, *The Annals of Statistics*, **37**, 3779–3821.
- [11] Muller, H. G. and Stadtmüller, U. (2005). Generalized functional linear models, *The Annals of Statistics*, **33**, 774–805.
- [12] Nelder, J. A. and Wedderburn, R. W. M. (1972). Generalized Linear Models, *Journal of Royal Statistical Society*, **135**, 370–384.
- [13] Ramsay, J. O. and Silverman, B. W. (1997). *Functional Data Analysis*. Springer, New York.

- [14] Rice, J. A. and Silverman, B. W. (1991). Estimating the mean and covariance structure nonparametrically when the data are curves. *Journal of the Royal Statistical Society*, **53**, 233–243.
- [15] Wedderburn, R. W. M. (1974). Quasi-likelihood functions, generalized linear models and the Gauss–Newton method. *Biometrika*, **61**, 439–447.