

## رگرسیون ناپارامتری

حمزه ترابی<sup>۱</sup>، مژگان بقایی پور<sup>۲</sup>

چکیده:

در این مقاله، پس از یادآوری رگرسیون پارامتری به بررسی معایب آن می‌پردازیم. سپس معرفی اجمالی بر رگرسیون ناپارامتری خواهیم داشت و در انتها دو روش از روش‌های رگرسیونی ساده‌ی ناپارامتری به نام‌های رگرسیون نادارایا - واتسون و رگرسیون چندجمله‌ای موضعی را معرفی خواهیم کرد. **واژه‌های کلیدی:** اعتبارسنجی متقابل، رگرسیون جمعی، رگرسیون چندجمله‌ای موضعی، رگرسیون ساده ناپارامتری.

### ۱ مقدمه

تحلیل رگرسیون یکی از شاخه‌های مهم علم آمار است که هدف آن، بررسی اثر متغیرهای مستقل (پیش‌بین) روی متغیر وابسته (پاسخ) است.

$$Y_i = m(x'_i) + \varepsilon_i,$$

معین می‌شود:

که در آن  $m(\cdot)$  تابع نامعلومی است و در اکثر روش‌های ناپارامتری، تابعی پیوسته و هموار (مشتق‌پذیر) در نظر گرفته می‌شود.

رگرسیون را از لحاظ نوع تابعی که رابطه‌ی بین متغیرهای مستقل و وابسته را تعیین می‌کند، می‌توانیم به دو دسته‌ی رگرسیون پارامتری و ناپارامتری تقسیم کنیم. در رگرسیون پارامتری، رابطه‌ی بین متغیرهای مستقل و وابسته از طریق تابع پارامتری  $f$  به صورت زیر معین می‌شود:

$$Y_i = f(\underline{\beta}', \mathbf{x}'_i) + \varepsilon_i,$$

در این مقاله به معرفی رگرسیون ناپارامتری پرداخته می‌شود. در ابتدا، برای فهم بهتر مطالب مربوط به رگرسیون ناپارامتری، به یادآوری مختصری از آن چه در دوره‌ی کارشناسی به عنوان رگرسیون پارامتری (سننتی) با آن آشنا شده‌ایم و بررسی معایب این روش می‌پردازیم [۱، ۲، ۸]. در ادامه، پس از معرفی کلی رگرسیون ناپارامتری به بحث و بررسی بیشتر درباره‌ی دو مدل از مدل‌های رگرسیونی ساده‌ی ناپارامتری به نام‌های رگرسیون نادارایا - واتسون و رگرسیون چندجمله‌ای موضعی پرداخته می‌شود. در این قسمت با ارائه‌ی

که در آن  $\varepsilon_i \sim (0, \Sigma^2)$  و  $\mathbf{x}'_i = (x_{i1}, \dots, x_{ip})$  و  $\underline{\beta}' = (\beta_1, \dots, \beta_p)$  برداری از پارامترهای مجهول است. در صورت برآورد  $\underline{\beta}$ ، تابع پارامتری  $f$  و بنابراین مدل به صورت  $\hat{y} = f(\underline{\hat{\beta}}', \mathbf{x}')$  برآورد می‌شود.

اما در رگرسیون ناپارامتری، رابطه‌ی بین متغیرهای مستقل و وابسته از طریق تابع ناپارامتری  $m(\cdot)$  به صورت زیر

<sup>۱</sup> عضو هیأت علمی دانشگاه یزد  
<sup>۲</sup> فارغ التحصیل کارشناسی ارشد از دانشگاه یزد

است. با به کارگیری نمادهای ماتریسی که در ادامه معرفی می‌شوند، مدل رگرسیونی چندجمله‌ای به صورت زیر بازنویسی می‌شود:

$$\mathbf{Y} = \mathbf{X}\underline{\beta} + \underline{\varepsilon},$$

که در آن نمادهای ماتریسی مذکور عبارت اند از:

$$\mathbf{X} = \begin{pmatrix} 1 & x_1 & x_1^2 & \dots & x_1^q \\ 1 & x_2 & x_2^2 & \dots & x_2^q \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & x_n^2 & \dots & x_n^q \end{pmatrix},$$

$$\mathbf{Y} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}, \quad \underline{\beta} = \begin{pmatrix} \beta_0 \\ \vdots \\ \beta_q \end{pmatrix}, \quad \underline{\varepsilon} = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix}.$$

با مشتق‌گیری از  $Q$  برحسب  $\beta_j$  ها و به کارگیری نمادهای ماتریسی، معادلات نرمال به صورت زیر به دست می‌آیند:

$$\mathbf{X}'\mathbf{X}\underline{\beta} = \mathbf{X}'\mathbf{Y}.$$

بنابراین برآورد پارامترهای نامعلوم مدل به صورت یکتا، از رابطه‌ی زیر حاصل می‌شود:

$$\underline{\hat{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}$$

و در آخر مدل به صورت زیر برآورد می‌شود:

$$\begin{aligned} \hat{\mathbf{Y}} &= \mathbf{X}\underline{\hat{\beta}} \\ &= \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y} \\ &= \mathbf{H}\mathbf{Y}, \end{aligned}$$

به طوری که  $\hat{\mathbf{Y}} = (\hat{y}_1 \dots \hat{y}_n)$

ماتریس  $\mathbf{H} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ ،  $n \times n$ ، که در واقع ماتریس تصویر، روی فضای تولید شده توسط ستون‌های

مثال‌های مختلف سعی در فهم بهتر مفاهیم و مطالب بیان شده داریم.

## ۲ یادآوری رگرسیون پارامتری

رگرسیون پارامتری را می‌توان به دو دسته رگرسیون پارامتری تک‌متغیره و چندمتغیره تقسیم‌بندی کرد.

(۱) تک‌متغیره:

$$Y_i = f(x_i) + \varepsilon_i,$$

مانند رگرسیون چندجمله‌ای که در آن

$$f(x_i) = \beta_0 + \beta_1 x_i + \dots + \beta_q x_i^q.$$

(۲) چندمتغیره:

$$Y_i = f(x_{i1}, \dots, x_{ip}) + \varepsilon_i,$$

مانند رگرسیون خطی چندگانه که در آن

$$f(x_{i1}, \dots, x_{ip}) = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}.$$

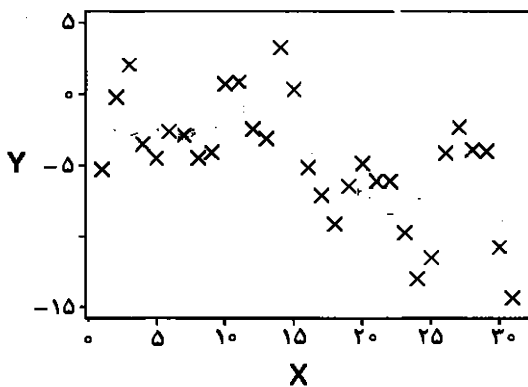
در این قسمت برای سادگی کار تنها مدل رگرسیونی تک‌متغیره چندجمله‌ای، یعنی مدلی با یک متغیر رگرسیونی  $x$  که با یک متغیر پاسخ  $Y$  دارای ارتباط است را مورد بررسی قرار می‌دهیم. داده‌ها به صورت  $(x_i, y_i)$ ،  $1 \leq i \leq n$  مدل

$$y_i = \sum_{j=0}^q \beta_j x_i^j + \varepsilon_i.$$

را در نظر بگیرید. در واقع هدف، برآورد پارامترهای نامعلوم  $\{\beta_0, \dots, \beta_q\}$  به روش کمترین مربعات خطا از طریق مینیمم کردن رابطه‌ی

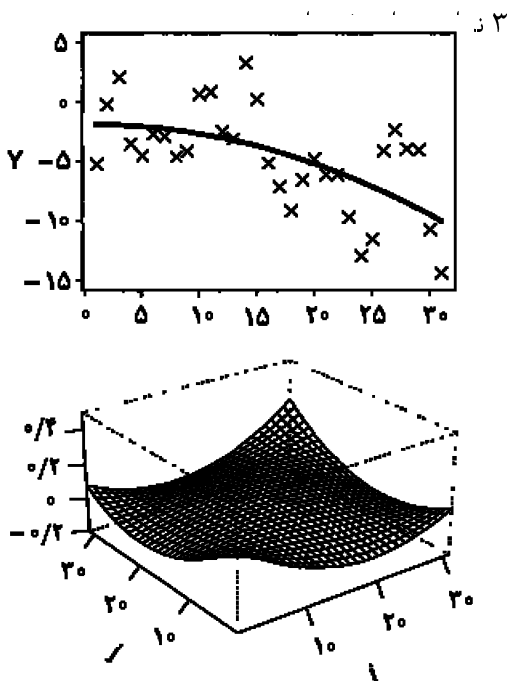
$$Q = \sum_{i=1}^n \left( \sum_{j=0}^q \beta_j x_i^j - y_i \right)^2.$$

مثال ۱-۲ میانگین روزانه‌ی درجه حرارت هوا در شهر واکانای<sup>۳</sup> (Y) برحسب روزهای ماه ژانویه‌ی ۱۹۹۶ ثبت شده است. روشن است که مقادیر متغیر مستقل، اعداد ۱، ۲، ۳، ...، ۳۰ هستند. نمودار پراکنش حاصل در شکل زیر نمایش داده شده است:



شکل ۱. نمودار پراکنش.

یک معادله‌ی چندجمله‌ای درجه‌ی ۲ در شکل ۲ و مقادیر درایه‌های ماتریس هت متناظر با برازش در شکل



ماتریس X است، ماتریس هت نام دارد که معمولاً ماتریس برازش نیز نامیده می‌شود. اگر ارتباط بین مقادیر مشاهدات  $y_j$  و مقادیر برازش شده‌ی  $\hat{y}_i$  را از طریق رابطه‌ی  $\hat{y}_i = \sum_{j=1}^n h_{ij} y_j$ ، مورد توجه قرار دهیم، در واقع می‌توانیم  $h_{ij}$  را به عنوان  $ij$  امین عضو ماتریس H در نظر بگیریم. ماتریس هت و ویژگی‌ها آن نقشی مهم در تحلیل رگرسیونی بازی می‌کنند.

به اختلاف بین مقادیر مشاهده شده‌ی  $y_i$  و مقادیر برازش شده‌ی  $\hat{y}_i$  متناظر  $\hat{y}_i$ ، یعنی  $e_i = y_i - \hat{y}_i$ ، باقیمانده می‌گوییم. مناسب است که  $n$  باقیمانده‌ی مذکور به صورت نماد ماتریسی زیر نوشته شود:

$$e = Y - \hat{Y},$$

هم چنین مجموع مربعات مانده‌ها عبارت است از

$$\begin{aligned} SSE &= e'e \\ &= Y'[I - H]Y. \end{aligned} \quad (1)$$

توجه کنید که در رگرسیون کاربردی، معادله‌ی رگرسیونی برازش شده‌ای مفید و مورد توجه است که دارای کمترین مقدار SSE باشد. همان طور که از رابطه‌ی ((۱)) مشخص است، کاهش SSE در رابطه‌ی مستقیم با ماتریس هت است به این صورت که اگر درایه‌های روی قطر اصلی ماتریس هت به یک و درایه‌های غیر قطر اصلی به صفر نزدیک باشند، میزان SSE کاهش می‌یابد. از این جاست که اهمیت ماتریس هت در مباحث رگرسیونی مشخص می‌شود. مثال زیر این نکته را روشن تر می‌کند.

داده‌ها یک ماکسیمم نسبی را نشان می‌دهند در حالی که نمودار برازش داده شده حاکی از یک مینیمم نسبی است. بنابراین، این نمودار برازش داده شده، به ویژه در مواقعی که هدف پیش‌بینی است، مناسب نخواهد بود. این نکته با توجه به نمودار ماتریس هت متناظر در شکل ۷ نیز نمایان است. این نمودار از نظر شکل ظاهری نوسانات زیادی دارد؛ یعنی ماتریس هت علاوه بر گرفتن مقادیر منفی، دارای مقادیر کوچکی از نظر قدرمطلق هم نیست و این مطلب با توجه به رابطه ((۱)) باعث افزایش چشم‌گیر SSE می‌شود. بنابراین همان طور که دیدید برازش یک معادله‌ی درجه ۷ نیز چندان مناسب به نظر نمی‌رسد.

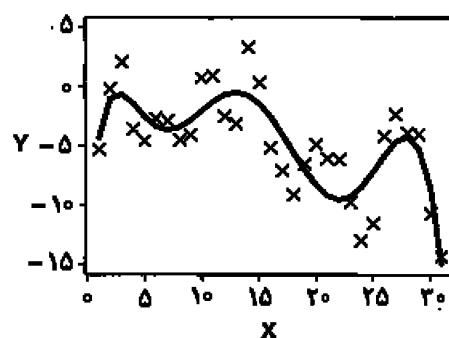
در واقع این نقص جزء ویژگی‌های ذاتی ماتریس هت در روش رگرسیون پارامتری چندجمله‌ای است که با افزایش درجه‌ی چندجمله‌ای مقادیر درایه‌های ماتریس هت به شدت افزایش پیدا می‌کنند.

اگرچه ماتریس هت یک آماره‌ی مهم را برای تشخیص رابطه‌ی بین داده‌ها و برآوردها فراهم می‌کند ولی این ماتریس همه‌ی ویژگی‌های رگرسیونی را پوشش نمی‌دهد. بنابراین برازش رگرسیون چندجمله‌ای که توسط ماتریس هت ارائه می‌شود، غالباً منجر به برآوردهایی دور از واقعیت می‌گردد و به کارگیری آن توصیه نمی‌شود.

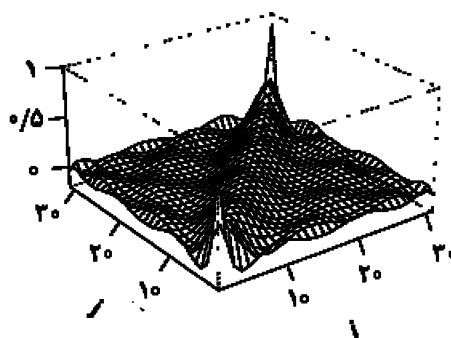
بنابراین باید از ماتریس هتی استفاده کنیم که یک برازش طبیعی را به داده‌ها ارائه دهد و این یکی از دلایلی است که از روش‌های بنا شده بر پایه‌ی رگرسیون ناپارامتری برای بررسی رابطه‌ی بین متغیرهای مستقل و پیش‌بین بهره می‌گیریم.

شکل ۳. درایه‌های ماتریس هت متناظر با شکل ۲.

همان طور که دیده می‌شود نمودارهای برازش داده شده از متغیرهای  $Y_i$  اختلاف زیادی دارند و این از معایب برازش معادله‌ی چندجمله‌ای محسوب می‌شود؛ برای نزدیک کردن مقادیر  $\hat{y}_i$  به متغیرهای  $Y_i$  از چندجمله‌ای‌هایی با درجه‌ی بالاتر بهره می‌گیریم. در ادامه یک معادله‌ی درجه‌ی ۷ در شکل ۴ به داده‌ها برازش می‌دهیم.



شکل ۴. برازش معادله‌ی درجه‌ی ۷.



شکل ۵. درایه‌های ماتریس هت متناظر با شکل ۶.

اگرچه در این نمودار مقادیر  $\hat{y}_i$  به داده‌های  $y_i$  نزدیک‌تر شده است اما باز هم نمودار برازش داده شده نمودار مناسبی به نظر نمی‌رسد. اولین مسأله‌ای که جلب توجه می‌کند در نقطه  $x_i = 22$  رخ داده است. در این نقطه

## ۳ معرفی رگرسیون ناپارامتری

یک نمودار به داده‌ها، منظور از هموارسازی، برازش یک نمودار مشتق‌پذیر به داده‌هاست.

$m(x)$  در معادله‌ی رگرسیون ساده‌ی ناپارامتری تابعی است که رفتار ذاتی داده‌ها را ارائه می‌دهد و به تصادف یا شانس وابسته نیست. مقدار برآورد شده‌ی  $m(x)$  را با  $\hat{m}(x)$  نمایش می‌دهیم که تابعی هموار از  $x$  است.

همان طور که در رگرسیون پارامتری گفتیم همیشه برازش چند جمله‌ای به روش کمترین مربعات مناسب نیست زیرا در این روش تنها یک چند جمله‌ای به ناحیه‌ای که داده‌ها در درون آن قرار گرفته‌اند، برازش داده می‌شود. برازش تنها یک معادله به تعداد زیادی از داده‌ها با رفتار پیچیده گاهی به نتایج نامطلوبی منجر می‌شود.

این مشکلات به مفهوم افراز کردن داده‌ها به چند ناحیه برای برازش منحنی به داده‌ها در هر ناحیه و عطف آن‌ها به طوری که برآوردگرهای مناسب و مفیدی باشند، منتهی می‌شود. در واقع ساده‌ترین روش، برازش معادله‌ی یک منحنی به صورت پاره‌ای (تکه‌ای) است به این شرط که برآوردها در انتهای هر زیر ناحیه مقادیر یکسانی را اختیار کنند.

انواع روش‌هایی که در رگرسیون ساده ناپارامتری برای برآورد  $m(x)$  وجود دارند، عبارتند از:

(۱) رگرسیون نادارایا - واتسون (برآوردگر هسته‌ای)

(۲) رگرسیون چند جمله‌ای موضعی

(۳) اسپلاین هموار ساز

(۴) روش LOESS

هدف از رگرسیون ناپارامتری، برآورد تابع رگرسیون به صورت مستقیم است. رگرسیون ناپارامتری اثر یک یا چند متغیر مستقل را روی یک متغیر وابسته بررسی می‌کند، بدون این که از قبل تابع ویژه‌ای را برای ارتباط بین متغیر مستقل و وابسته در نظر گرفته باشد. رگرسیون ناپارامتری را می‌توان از لحاظ تعداد متغیرهای مستقل به دو گروه زیر تقسیم‌بندی کرد:

(۱) تک‌متغیره (رگرسیون ساده‌ی ناپارامتری):

$$Y_i = m(x_i) + \varepsilon_i.$$

(۲) چندمتغیره (رگرسیون جمعی ناپارامتری):

$$Y_i = m_0 + m_1(x_{i1}) + m_2(x_{i2}) + \dots + m_k(x_{ik}) + \varepsilon_i, \quad (2)$$

به طوری که به  $m_j(\cdot)$  ها،  $j = 1, \dots, k$ ، تابع‌های

رگرسیونی جزئی گفته می‌شود و  $\varepsilon_i \sim (0, \Sigma^2)$ .

رگرسیون ساده‌ی ناپارامتری یکی از ساده‌ترین مدل‌های رگرسیونی ناپارامتری است که در آن متغیر مستقل به صورت تک‌متغیره در نظر گرفته می‌شود. در این مدل، داده‌ها به صورت  $(x_i, y_i)$ ،  $i = 1, 2, \dots, n$ ، هستند که غالباً  $x_i$  ها به صورت صعودی مرتب شده‌اند و داریم:

$$m(x_i) = E(Y_i | X_i = x_i).$$

با توجه به این که در این روش، نمودار پراکنش  $y_i$  بر حسب  $x_i$  کاربرد زیادی دارد، در برخی از مراجع نام هموارسازی پراکنش نیز برای این روش در نظر گرفته شده است. توجه کنید که در انواع روش‌های برازش

تیبشیرانی<sup>۶</sup> [۴] معرفی و گسترش داده شد، مشهورترین روش برای برآورد مدل جمعی بوده است. این روش در مطالعات شبیه‌سازی و در به کارگیری داده‌های واقعی روشی مناسب به حساب می‌آید اما به علت طبیعت تکراری این روش، تحلیل ویژگی‌های آماری مدل با مشکلات زیادی روبه‌رو بود. به همین دلیل روش انتگرال‌گیری کناری توسط تی‌جوستیم<sup>۷</sup> و اوستاد<sup>۸</sup> در سال ۱۹۹۴ ولین‌تون و نیلسن<sup>۹</sup> [۵] معرفی شد. برای آشنایی بیشتر با این دو روش می‌توانید به مراجع [۴]، [۵] و [۶] مراجعه کنید.

در ادامه تنها با دو روش اول از رگرسیون ساده‌ی ناپارامتری آشنا می‌شویم. برای آشنایی بیشتر با رگرسیون ناپارامتری می‌توانید به فاکس<sup>۱۰</sup> [۳]، سیمنوف<sup>۱۱</sup> [۷] و تاکزاوا<sup>۱۲</sup> [۸] مراجعه کنید.

### ۱.۳ برآوردگر نادارایا - واتسون

در ادامه به معرفی روش رگرسیون برآوردگر نادارایا - واتسون می‌پردازیم که از میانگین وزنی داده‌ها به طور مستقیم برای برآورد  $m(x)$  بهره می‌گیرد.

فرض کنید که متغیرهای مستقل و وابسته، هر دو متغیر تصادفی پیوسته باشند. می‌دانیم که هرگاه  $f_X(x)$  تابع چگالی  $X$  باشد آنگاه برای هر  $\Delta x$  کوچک

$$P(x < X < x + \Delta x) \approx f_X(x)\Delta x, \quad (۴)$$

اگرچه در نگاه اول ممکن است این نکته مشهود نباشد، اما تمام روش‌هایی که در این قسمت با آن‌ها آشنا می‌شویم در چارچوب کلی میانگین وزنی جای می‌گیرند.

اما در رگرسیون جمعی که توسط رابطه‌ی ((۲)) معرفی شد، داده‌های نمونه‌ی تصادفی به صورت  $\{(x_i, y_i)\}_{i=1}^n$  هستند، به طوری که  $x_i \in \mathbb{R}^k$  و  $y_i \in \mathbb{R}$ . در رابطه‌ی ((۲))،  $m_0$  مقداری ثابت است که توسط  $E(Y)$  برآورد می‌شود و تابع‌های رگرسیونی جزئی  $n$ ،  $j = 1, 2, \dots, n$  که به وسیله‌ی داده‌ها برآورد می‌شوند، تابع‌هایی هموار هستند به طوری که برای هر  $j > 0$

$$E[m_j(X_j)] = 0.$$

این مدل را می‌توان در دو جهت مختلف گسترش داد:

(۱) برای اضافه کردن اثر متقابل بین متغیرهای مستقل مورد نظر، مانند

$$Y_i = m_0 + m_1(x_{1i}) + m_{2,3}(x_{2i}, x_{3i}).$$

(۲) برای اضافه کردن جملات خطی در مدل، مانند

$$Y_i = m_0 + \beta_1 x_{1i} + m_2(x_{2i}), \quad (۳)$$

مدل ارائه شده در ((۳)) مدل نیمه‌پارامتری نامیده می‌شود.

از حدود پانزده سال قبل روش قدیمی پس‌برازشی که توسط بیوجا<sup>۴</sup> و همکاران در سال ۱۹۸۹ و هاستای<sup>۵</sup> و

Buja<sup>۴</sup>Hastie<sup>۵</sup>Tibshirani<sup>۶</sup>Tjøstheim<sup>۷</sup>Auestad<sup>۸</sup>Linton and Nielsen<sup>۹</sup>Fox<sup>۱۰</sup>Simonoff<sup>۱۱</sup>Takezawa<sup>۱۲</sup>

تابع‌های هسته‌ای معمولاً در شرایط زیر صدق می‌کنند:

$$\int_{-\infty}^{\infty} k_x(u) du = \int_{-\infty}^{\infty} k_y(u) du = 1, \quad (۷)$$

$$\int_{-\infty}^{\infty} uk_x(u) du = \int_{-\infty}^{\infty} uk_y(u) du = 0, \quad (۸)$$

$$\int_{-\infty}^{\infty} u^2 k_x(u) du < \infty,$$

$$\int_{-\infty}^{\infty} u^2 k_y(u) du < \infty.$$

از معادلات بیان شده در روابط ((۷)) و ((۸)) نتیجه می‌گیریم که

$$\frac{1}{h_x} \int_{-\infty}^{\infty} k_x\left(\frac{x - X_i}{h_x}\right) dx = \frac{1}{h_y} \int_{-\infty}^{\infty} k_y\left(\frac{y - Y_i}{h_y}\right) dy = 1,$$

$$\frac{1}{h_x} \int_{-\infty}^{\infty} x k_x\left(\frac{x - X_i}{h_x}\right) dx = X_i,$$

$$\frac{1}{h_y} \int_{-\infty}^{\infty} y k_y\left(\frac{y - Y_i}{h_y}\right) dy = Y_i. \quad (۹)$$

به این ترتیب برای به دست آوردن  $\hat{m}(x)$  کافیست برآوردهای  $f_X(x)$  و  $f(x, y)$  را در رابطه‌ی ((۶)) جایگزین کنیم بنابراین

$$\begin{aligned} \hat{m}(x) &= \int_{-\infty}^{\infty} y \frac{\hat{f}(x, y)}{\hat{f}_X(x)} dy \\ &= \int_{-\infty}^{\infty} y \frac{\frac{1}{nh_x h_y} \sum_{i=1}^n k_x\left(\frac{x - X_i}{h_x}\right) k_y\left(\frac{y - Y_i}{h_y}\right)}{\frac{1}{nh_x} \sum_{i=1}^n k_x\left(\frac{x - X_i}{h_x}\right)} dy \\ &= \frac{\sum_{i=1}^n k_x\left(\frac{x - X_i}{h_x}\right) \int_{-\infty}^{\infty} y k_y\left(\frac{y - Y_i}{h_y}\right) dy}{h_y \sum_{i=1}^n k_x\left(\frac{x - X_i}{h_x}\right)}. \quad (۱۰) \end{aligned}$$

با جایگزین کردن رابطه‌ی (۹) در (۱۰) و جایگزین کردن  $k_x(\cdot)$  و  $h_x$  با  $k(\cdot)$  و  $h$  به دست می‌آوریم:

$$\hat{m}(x) = \sum_{i=1}^n w_i(x, h) Y_i, \quad (۱۱)$$

به همین ترتیب هرگاه  $f(x, y)$  تابع چگالی توأم  $X$  و  $Y$  باشد آنگاه برای هر  $\Delta x$  و  $\Delta y$  کوچک

$$\begin{aligned} P(x < X < x + \Delta x, y < Y < y + \Delta y) \\ \approx f(x, y) \Delta x \Delta y, \quad (۵) \end{aligned}$$

هم‌چنین هرگاه  $f(y|x)$  تابع چگالی احتمال شرطی  $Y$  به شرط  $x$  باشد آنگاه

$$P(y < Y < y + \Delta y | x < X < x + \Delta x) \approx f(y|x) \Delta y,$$

با توجه به این نکته که مقدار تابع  $m(x)$  برای مقدار مستقل  $x$  هم‌ارز با امید ریاضی متغیر  $Y$  به شرط  $X = x$  است، به دست می‌آوریم که

$$\begin{aligned} m(x) &= E(Y | X = x) \\ &= \int_{-\infty}^{\infty} y f(y|x) dy \\ &= \frac{1}{f_X(x)} \int_{-\infty}^{\infty} y f(x, y) dy. \quad (۶) \end{aligned}$$

رابطه‌ی ((۶)) نشان می‌دهد که برای برآورد تابع  $m(x)$  کافی است تابع  $f(y|x)$  یا در واقع توابع  $f_X(x)$  و  $f(x, y)$  را برآورد کنیم. یکی از ساده‌ترین روش‌ها برای برآورد توابع  $f_X(x)$  و  $f(x, y)$  به کارگیری برآوردگر چگالی هسته‌ای است. با استفاده از ایده‌ی روش گشتاورها و روابط ((۴)) و ((۵))، برآوردهای توابع  $f_X(x)$  و  $f(x, y)$  به صورت زیر به دست می‌آیند:

$$\hat{f}(x, y) = \frac{1}{nh_x h_y} \sum_{i=1}^n k_x\left(\frac{x - X_i}{h_x}\right) k_y\left(\frac{y - Y_i}{h_y}\right),$$

$$\hat{f}_X(x) = \frac{1}{nh_x} \sum_{i=1}^n k_x\left(\frac{x - X_i}{h_x}\right),$$

به طوری که  $k_x(\cdot)$  و  $k_y(\cdot)$  تابع‌های هسته‌ای  $h_x$  و  $h_y$  مقادیری ثابت و مثبت به نام پهنه هستند که در ادامه بیشتر با آن‌ها آشنا می‌شویم.

به طوری که

$$w_i(x, h) = \frac{k\left(\frac{x-X_i}{h}\right)}{\sum_{k=1}^n k\left(\frac{x-X_k}{h}\right)}$$

برآوردگری که برای برآورد  $m(x)$ ، در رابطه‌ی (۱۱) معرفی شد، برآوردگر نادارایا - واتسون نام دارد.

تابع  $w_i(x, h)$  که هسته هم‌ارز نامیده می‌شود، وزن‌هایی را معین می‌کند که برای ارزیابی چگونگی اهمیت  $x_i$  ها در رابطه با  $Y_i$ ، مورد توجه قرار می‌گیرد. با توجه به رابطه‌ی (۱۱) دیده می‌شود که این برآوردگر، یک برآوردگر خطی است. نکته قابل توجه دیگر این که  $w_i(x, h)$  در برآوردگر نادارایا - واتسون دارای ویژگی زیر است:

$$\sum_{i=1}^n w_i(x, h) = 1.$$

یعنی برآوردگر نادارایا - واتسون میانگین وزنی  $Y_i$  ها است. روشن است که تابع  $w_i(x, h)$  در رابطه‌ی مستقیم با مقدار  $h$  و تابع  $k(\cdot)$  است.

توجه کنید که شرط بیان شده در رابطه‌ی (۷) برای توابع هسته‌ای  $k(\cdot)$  قابل صرف نظر کردن است یعنی اگر تابع هسته‌ای در این شرط صدق نکند، در به دست آوردن  $\hat{m}(x)$  تأثیری ندارد. یکی از ساده‌ترین تابع‌های هسته‌ای، تابع  $\exp(-u^2)$  است که به صورت‌های زیر نیز می‌توان آن را به کار گرفت:

$$\exp\left(\frac{-u^2}{2}\right), \quad \frac{1}{\sqrt{2\pi}} \exp\left(\frac{-u^2}{2}\right).$$

بهره‌گیری از هر کدام از توابع بالا، برآورد تابع رگرسیونی  $\hat{m}(x)$  را بدون تغییر باقی می‌گذارد.

پهنه‌ی  $h$  در رابطه‌ی (۱۱) را می‌توان به صورت دنباله‌ای مثبت  $\{h_n\}$  در نظر گرفت که با افزایش حجم نمونه به

سمت صفر میل می‌کند. از  $h$  می‌توان به عنوان پارامتر هموارسازی نیز یاد کرد؛ به این مفهوم که هرچه مقدار  $h$  افزایش (کاهش) یابد، نمودار حاصل از برازش رگرسیون ناپارامتری هموارتر (ناهموارتر) می‌گردد. نکته‌ای که در این جا باید به آن توجه کنیم این است که افزایش  $h$  نمی‌تواند به صورت دلخواه ادامه یابد زیرا افزایش  $h$  در رابطه‌ی مستقیم با افزایش مربع انحرافات از میانگین<sup>۱۳</sup> است و همان طور که می‌دانیم برآورد یک تابع رگرسیونی با MSE زیاد مطلوب نیست؛ بنابراین  $h$  را باید به صورتی انتخاب کنیم که نمودار حاصل هموارترین نمودار با کوچک‌ترین MSE باشد. در ادامه دو معیار برای انتخاب  $h$  بهینه معرفی خواهیم کرد.

برآوردگر  $Y_i$ ، یعنی  $\hat{Y}_i$  را می‌توان به صورت زیر مشخص نمود:

$$\hat{Y}_i = \sum_{j=1}^n h_{ij} Y_j,$$

که در آن

$$h_{ij} = w_j(x_i, h) = \frac{k\left(\frac{x_i - x_j}{h}\right)}{\sum_{k=1}^n k\left(\frac{x_i - x_k}{h}\right)}, \quad (12)$$

که  $h_{ij}$  در رابطه‌ی (۱۲)،  $ij$  امین عضو ماتریس  $H$  خواهد بود. اعضای قطری ماتریس هتی که در رابطه‌ی (۱۲) معرفی شد، کاربردهای مفیدی دارند که در ادامه به بررسی آن‌ها می‌پردازیم.

### ۲.۳ شاخصی سودمند برای انتخاب معادله‌ی رگرسیونی

همان طور که ملاحظه کردید، تابع  $\hat{m}(x)$  از طریق تابع وزنی  $w_i(x, h)$ ، در ارتباط با پهنه‌ی  $h$  است و با تغییر

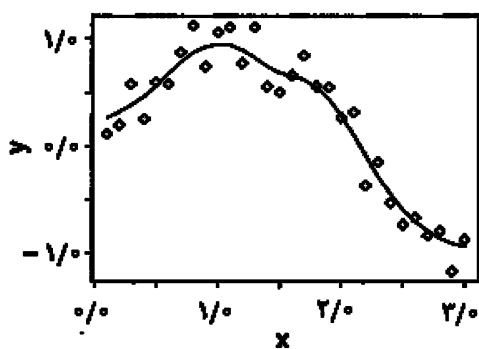


تعریف ۳-۲ شاخص اعتباربخشی متقابل تعمیم یافته، شاخصی برای مقایسه‌ی بین نمودارهای رگرسیونی است که به صورت زیر تعریف می‌شود:

$$GCV[\hat{m}(x)] = \frac{\sum_{i=1}^n (Y_i - \hat{m}(X_i))^2}{n(1 - \frac{\sum_{i=1}^n h_{ii}}{n})^2}$$

روشن است که  $h$  بهینه مربوط به مینیمم CV و GCV است. در این جا برای آشنایی بیشتر با این روش مثالی ارائه می‌کنیم.

مثال ۳-۲ در شکل ۶ نموداری با استفاده از روش رگرسیونی برآوردگر نادارایا - واتسون به داده‌هایی که با علامت  $\diamond$  مشخص شده‌اند، برازش داده‌ایم. در این برازش  $h$  بهینه برابر با  $0.18$  در نظر گرفته شده است. این انتخاب با استفاده از معیار CV که در شکل ۷ نمایش داده شده است، صورت گرفته است.



شکل ۶. برازش نمودار به وسیله‌ی رگرسیون نادارایا - واتسون.

مقدار  $h$  می‌توان هموارسازی‌های متفاوتی را از  $\hat{m}(x)$  به دست آورد. تعیین مقدار پهنه به نحوی که نمودار رگرسیونی برازش داده شده نموداری مفید باشد، نکته‌ای قابل توجه است. بنابراین باید از بین نمودارهایی که می‌توان به داده‌ها برازش داد، بهترین را انتخاب نمود. اما رسم همه‌ی نمودارهای ممکن و سپس انتخاب یکی از بین آن‌ها کار عاقلانه‌ای به نظر نمی‌رسد و غیر ممکن است. به همین دلیل معرفی شاخص‌هایی برای انتخاب بهترین مدل از طریق انتخاب  $h$  الزامی به نظر می‌رسد.

یکی از شاخص‌های سودمند و کاربردی برای مقایسه‌ی نمودارهای رگرسیونی، شاخص اعتبارسنجی متقابل است.

تعریف ۳-۱ شاخص اعتبارسنجی متقابل  $\hat{m}(x)$  شاخصی برای مقایسه‌ی پیش‌بینی داده‌های مربوط به آینده است و فرمول آن از رابطه‌ی زیر به دست می‌آید:

$$CV[\hat{m}(x)] = \sum_{i=1}^n \frac{(Y_i - \hat{m}(X_i))^2}{n(1 - h_{ii})^2}$$

نکته‌ای که باید در بهره‌گیری از شاخص CV مورد توجه قرار دهیم این است که این رابطه برای معادله‌های رگرسیونی که در شرایط زیر صدق می‌کنند، مفید است:

(۱) معادله‌ی رگرسیونی به روش کمترین مربعات به وسیله‌ی ماتریس هت به دست آمده باشد.

(۲) معادله‌ی رگرسیونی مورد استفاده، یکتا باشد.

شاخص دیگری که برای مقایسه‌ی معادله‌های رگرسیونی به کار گرفته می‌شود، شاخص اعتباربخشی متقابل تعمیم یافته نام دارد.

### ۴.۳ رگرسیون چندجمله‌ای موضعی

رگرسیون چندجمله‌ای موضعی یکی از انواع روش‌های رگرسیونی ساده‌ی ناپارامتری است. در این روش نیز مانند اکثر روش‌های رگرسیون ناپارامتری ناحیه‌ای را که داده‌ها در آن قرار گرفته‌اند به زیرناحیه‌هایی تقسیم می‌کنیم و سپس معادله‌ی خطی را به هر زیرناحیه برآزش می‌دهیم به طوری که نمودارها در انتهای هر زیرناحیه پیوسته باشند.

همان طور که گفته شد رگرسیون چندجمله‌ای موضعی، تعمیمی از روش برآوردگر نادارایا - واتسون است. معادله‌ی چندجمله‌ای موضعی وقتی متغیر مستقل  $x$  به  $x^*$  نزدیک است به صورت زیر تعریف می‌شود:

$$m(x, x^*) = a_0(x^*) + \sum_{j=1}^p a_j(x^*)(x - x^*)^j. \quad (13)$$

ساده‌ترین و عمومی‌ترین حالت در این روش، حالتی است که  $p = 1$  در نظر گرفته شود؛ یعنی از روش رگرسیون خطی موضعی بهره بگیریم.

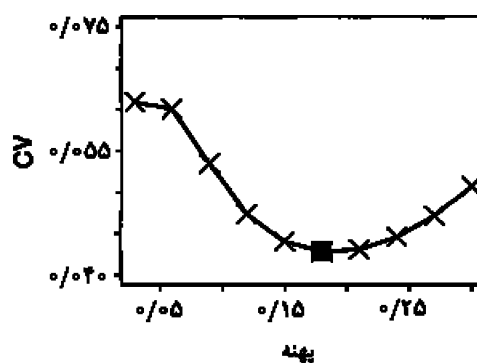
مقادیر  $\{a_0(x^*), a_1(x^*), \dots, a_p(x^*)\}$  از مینیمم کردن رابطه‌ی زیر به دست می‌آیند:

$$E_{local}(x^*) = \sum_{i=1}^n \left( w \left( \frac{x_i - x^*}{h} \right) (m(x_i, x^*) - y_i)^2 \right) \quad (14)$$

که در آن

$$m(x_i, x^*) = a_0(x^*) + \sum_{j=1}^p a_j(x^*)(x_i - x^*)^j.$$

$w \left( \frac{x_i - x^*}{h} \right)$  در رابطه‌ی (۱۴) همان تابع هسته‌ای است که درجه اهمیت نقطه‌ی  $(x_i, y_i)$ ، وقتی معادله‌ی رگرسیونی مربوط به  $x^*$  مدنظر است را معین می‌کند. این تابع وزنی نامنفی، بیشترین مقدار خود را در حالت  $x_i = x^*$



شکل ۷. نتایج محاسبه‌ی CV برای داده‌های شکل ۶.

### ۳.۳ اریبی کران

یک اشکال ذاتی در روش برآوردگر هسته‌ای وجود دارد. این اشکال وقتی رخ می‌دهد که  $x_i$  که می‌خواهیم  $\hat{y}_i$  متناظرش را برآورد کنیم یکی از کران‌های داده‌ها یعنی مینیمم یا ماکسیمم داده‌ها باشد. این محدودیت از چپ یا راست  $x_i$  باعث می‌شود که  $\hat{y}_i$  به شدت تحت تأثیر داده‌هایی قرار گیرد که  $x_i$  در آن سمت محدود نیست. در واقع در این موقعیت، پهنه‌ی  $h$  به طور متقارن گسترده نشده است و به مقادیری که در سمت کران  $x_i$  واقع شده‌اند، به ناچار وزن صفر داده می‌شود. این نکته باعث می‌شود که  $\hat{y}_i$  برآورد شده به شدت تحت تأثیر داده‌هایی که در جهت دیگر قرار گرفته‌اند، واقع شود. به این اشکال در رگرسیون نادارایا - واتسون اریبی کران گفته می‌شود.

رگرسیون چندجمله‌ای موضعی، یکی از انواع روش‌های رگرسیون ساده‌ی ناپارامتری است که این اشکال را به صورت چشم‌گیری تعدیل می‌کند و در واقع می‌توان روش رگرسیون چندجمله‌ای موضعی را به عنوان تعمیمی از رگرسیون نادارایا - واتسون در نظر گرفت.

می‌گیرد. انواع مختلف توابع وزنی که در این روش

می‌توانیم از آن‌ها بهره بگیریم عبارتند از:

$$w\left(\frac{x_i - x^*}{h}\right) = \exp\left(-\frac{1}{\nu}\left(\frac{x_i - x^*}{h}\right)^\nu\right),$$

$$= \begin{cases} \left(1 - \left(\frac{x_i - x^*}{h}\right)^\nu\right)^\nu & \text{if } \left(\frac{x_i - x^*}{h}\right)^\nu \leq 1 \\ 0 & \text{if } \left(\frac{x_i - x^*}{h}\right)^\nu > 1, \end{cases}$$

$$= \begin{cases} \left(1 - \left(\frac{x_i - x^*}{h}\right)^\nu\right)^\nu & \text{if } \left(\frac{x_i - x^*}{h}\right)^\nu \leq 1 \\ 0 & \text{if } \left(\frac{x_i - x^*}{h}\right)^\nu > 1, \end{cases}$$

$$= \begin{cases} \left(1 - \left|\frac{x_i - x^*}{h}\right|^\nu\right)^\nu & \text{if } \left|\frac{x_i - x^*}{h}\right|^\nu \leq 1 \\ 0 & \text{if } \left|\frac{x_i - x^*}{h}\right|^\nu > 1. \end{cases}$$

با بهره‌گیری از نمادهای ماتریسی،  $E_{local}(x^*)$  را به صورت ساده‌تر زیر می‌نویسیم:

$$E_{local}(x^*) = (\mathbf{X}\mathbf{a} - \mathbf{Y})' \mathbf{W} (\mathbf{X}\mathbf{a} - \mathbf{Y}), \quad (15)$$

به کمک رابطه (۱۵) برآوردگر  $\mathbf{a}$  از رابطه‌ی

$$\hat{\mathbf{a}} = (\mathbf{X}' \mathbf{W} \mathbf{X})^{-1} \mathbf{X}' \mathbf{W} \mathbf{Y},$$

به دست می‌آید که در آن ماتریس  $\mathbf{X}$ ، ماتریس وزنی  $\mathbf{W}$

و بردارهای  $\mathbf{y}$  و  $\mathbf{a}$  به صورت زیر تعریف می‌شوند:

$$\mathbf{X} = \begin{pmatrix} 1 & (x_1 - x^*) & \dots & (x_1 - x^*)^p \\ 1 & (x_2 - x^*) & \dots & (x_2 - x^*)^p \\ \vdots & \vdots & \ddots & \vdots \\ 1 & (x_n - x^*) & \dots & (x_n - x^*)^p \end{pmatrix},$$

$$\mathbf{W} = \begin{pmatrix} w\left(\frac{x_1 - x^*}{h}\right) & 0 & \dots & 0 \\ 0 & w\left(\frac{x_2 - x^*}{h}\right) & \dots & 0 \\ 0 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & w\left(\frac{x_n - x^*}{h}\right) \end{pmatrix},$$

نکته‌ای که باید در این جا به آن توجه کنیم این است که ماتریس  $\mathbf{X}$  که در این روش به کار می‌بریم، تابعی از  $x^*$  است. در واقع در رگرسیون چندجمله‌ای، معادله‌ی رگرسیونی به کل داده‌هایی که در مدل وجود دارند، برازش داده می‌شود؛ در حالی که در رگرسیون چندجمله‌ای موضعی، معادله‌ی رگرسیونی را به داده‌هایی که در همسایگی  $x^*$  وجود دارند، برازش می‌دهیم.

با توجه به رابطه‌ی (۱۳) مقدار  $\hat{m}(x^*, x^*)$  برابر با  $\hat{a}_0(x^*)$  است که با به کارگیری فرم ماتریسی می‌توانیم آن را به صورت زیر بنویسیم:

$$\hat{m}(x^*, x^*) = \mathbf{e}'_1 \left( (\mathbf{X}' \mathbf{W} \mathbf{X})^{-1} \mathbf{X}' \mathbf{W} \mathbf{Y} \right),$$

که در آن  $\mathbf{e}_1$  بردار ستونی  $(1, 0, 0, \dots, 0)'$  با  $p + 1$  عضو است. فرم دیگر نمایش  $\hat{m}(x^*, x^*)$  به صورت زیر است:

$$\hat{m}(x^*, x^*) = \mathbf{q}(x^*)' \mathbf{Y}$$

$$= \sum_{j=1}^n q_j(x^*) y_j, \quad (16)$$

که در آن  $\mathbf{q}(x^*)$  از رابطه‌ی زیر پیروی می‌کند:

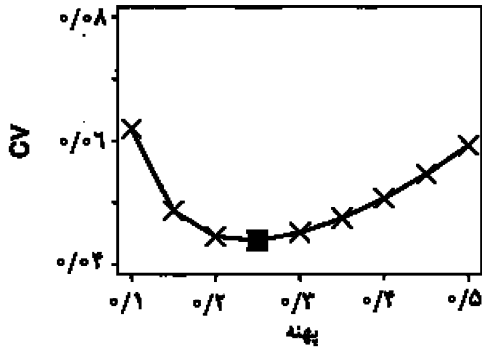
$$\mathbf{q}(x^*)' = \mathbf{e}'_1 (\mathbf{X}' \mathbf{W} \mathbf{X})^{-1} \mathbf{X}' \mathbf{W}. \quad (17)$$

حال دو طرف رابطه‌ی (۱۷) را از سمت راست بردار در  $\mathbf{X}$  ضرب می‌کنیم، داریم:

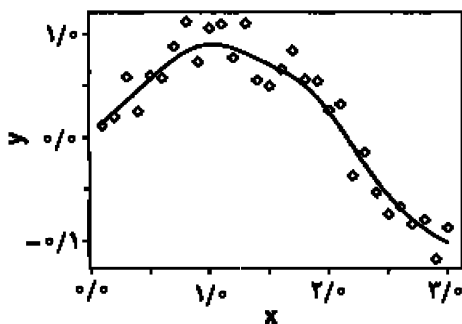
$$\mathbf{q}(x^*)' \mathbf{X} = \mathbf{e}'_1 (\mathbf{X}' \mathbf{W} \mathbf{X})^{-1} \mathbf{X}' \mathbf{W} \mathbf{X}$$

$$= \mathbf{e}'_1.$$

$$\mathbf{a} = \begin{pmatrix} a_0(x^*) \\ \vdots \\ a_p(x^*) \end{pmatrix}.$$



شکل ۸. نتایج محاسبه‌ی CV برای داده‌های شکل ۶ در روش رگرسیون خطی موضعی.



شکل ۹. برازش نمودار به وسیله‌ی روش رگرسیون خطی موضعی.

همان طور که می‌بینید، در شکل ۸ مقدار  $h$  بهینه با شاخص CV برابر با ۰/۲۵ تعیین شده است. با توجه به مقادیر  $h$  که در بالا معین شد و با بهره‌گیری از تابع وزنی چگالی نرمال، در شکل ۹ نموداری را به داده‌ها برازش می‌دهیم. برازش با خط پیوسته برای داده‌هایی که به وسیله‌ی  $\diamond$  مشخص شده‌اند، نمایش داده شده است.

از مقایسه‌ی شکل ۹ با شکل ۶ متوجه می‌شویم که نمودار برازش داده شده در شکل ۹ کمتر تحت تأثیر نقاطی که داده‌ها در آن کراندارند، قرار گرفته است که این ویژگی از جمله امتیازات روش رگرسیون چندجمله‌ای موضعی

با توجه به این نکته که  $x_{i1}$  ها،  $1 \leq i \leq n$  همگی برابر با یک هستند، رابطه‌ی زیر برقرار است:

$$\sum_{i=1}^n q_i(x^*) = 1. \quad (18)$$

با توجه به روابط (۱۶) و (۱۸) نتیجه می‌گیریم که برآورد حاصل از روش رگرسیون چندجمله‌ای موضعی، همان میانگین وزنی  $Y_j$  ها است. اگر  $x^* = x_k$  آنگاه

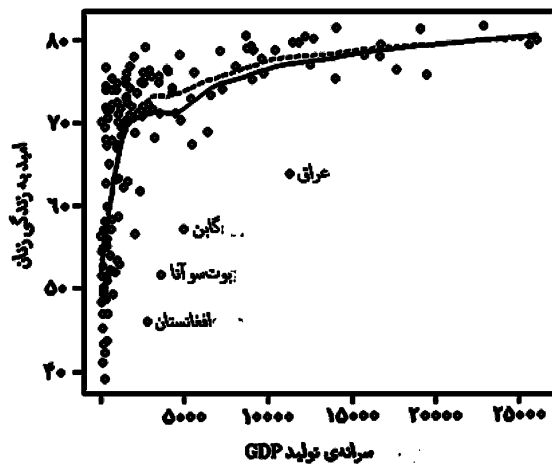
$$\begin{aligned} \hat{Y}_k &= \hat{m}(x_k, x_k) \\ &= \sum_{j=1}^n q_j(X_k) Y_j, \end{aligned}$$

که در آن  $q_j(x_k)$ ،  $k, j$  امین عضو ماتریس هت است. در این حالت  $k$  امین ستون ماتریس  $X'$  مانند بردار  $e'_k$  است و می‌توانیم  $k$  امین عضو روی قطر ماتریس هت را به صورت زیر تعریف کرد:

$$\begin{aligned} h_{kk} &= q_k(X_k) \\ &= [(X'WX)^{-1}]_{11} \cdot w(0). \end{aligned}$$

با توجه به ماتریس هت روش رگرسیون چندجمله‌ای موضعی، مقادیر CV و GCV را برای تعیین پهنه‌ی  $h$  بهینه به کار می‌بریم. برای بررسی بیشتر جزئیات ذکر شده در این قسمت به [۸] صفحات ۱۵۱-۱۴۱ مراجعه کنید.

مثال ۳-۳ در این جا به داده‌های مثال قبلی، نموداری را به وسیله‌ی روش رگرسیون خطی موضعی برازش می‌دهیم. ابتدا با کمک نمودار مربوط به CV مقدار  $h$  بهینه را انتخاب می‌کنیم. این نمودارها در شکل زیر نمایش داده شده‌اند.



شکل ۱۰. برازش نمودار به وسیله‌ی روش رگرسیون خطی موضعی.

حال می‌خواهیم به وسیله‌ی رگرسیون نادارایا - واتسون به صورت گام به گام، نموداری را به این داده‌ها برازش دهیم. (۱) در گام اول یکی از مقادیر  $x_i$  ها را انتخاب می‌کنیم؛ برای مثال می‌توانید  $x_i$  ها را به صورت صعودی مرتب کرده و سپس یکی یکی آن‌ها را انتخاب کنیم. به  $x_i$  که در این مرحله انتخاب می‌کنیم، اصطلاحاً مقدار  $x$  کانونی می‌گوییم و آن را با  $x_0$  (یا همان طور که در متن به کار گرفته شد  $x^*$ ) نمایش می‌دهیم. در این جا برای مثال  $x_{(۸۰)}$  را انتخاب می‌کنیم.

حال نوبت به انتخاب پهنه‌ی  $h$  می‌رسد. برای انتخاب  $h$  راه‌های مختلفی وجود دارد؛ برای مثال می‌توانید  $m$  امین همسایگی نزدیک  $x_0$  را در نظر بگیرید. در این جا  $h$  نصف طول بازه‌ی در برگرفته‌ی  $m$  امین همسایگی نزدیک  $x_0$  خواهد بود. این ناحیه در شکل ۱۱ (الف) مشخص شده است. دقت کنید که هر چه  $h$  بزرگ‌تر باشد، نمودار حاصل هموارتر خواهد شد.

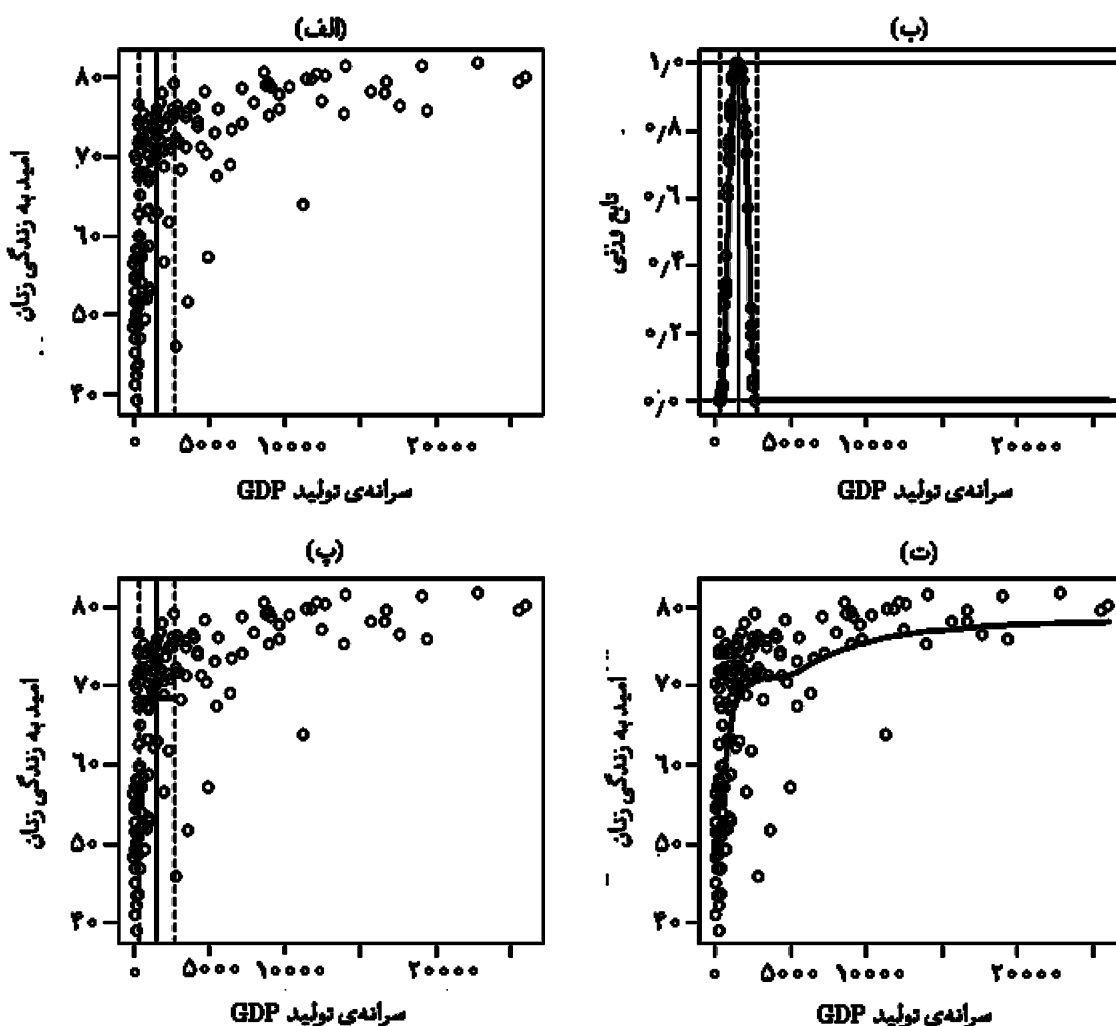
نسبت به رگرسیون نادارایا - واتسون است.

این مقاله را با ارائه‌ی مثال دیگری به پایان می‌رسانیم. در این مثال با بررسی گام به گام آن چه برای رسم یک نمودار در حالت رگرسیون ناپارامتری ساده انجام می‌دهیم، برآینم تا درک درستی از این روش در ذهن مخاطبین ایجاد کنیم. البته رسم نمودار با روش‌های ناپارامتری به صورت دستی کار عاقلانه‌ای نیست و در عمل غیر ممکن است و در این موارد از برنامه‌های کامپیوتری کمک می‌گیریم [۸].

مثال ۳-۴ در شکل ۱۰ داده‌های مربوط به امید زندگی زنان و سرانه‌ی تولید ناخالص داخلی<sup>۱۴</sup> (GDP) نمایش داده شده‌اند. هدف، بررسی ارتباط بین امید به زندگی زنان با سرانه‌ی تولید GDP است. داده‌های مذکور از ۱۵۴ کشور جهان در سال ۱۹۹۸ توسط سازمان ملل جمع‌آوری شده است. در این شکل دو نمودار را به وسیله‌ی روش رگرسیون خطی موضعی به داده‌ها برازش داده‌ایم. نمودار پیوسته برازش را روی همه‌ی داده‌ها انجام داده است در حالی که نمودار نقطه چین چهار کشوری را که روی شکل مشخص شده‌اند، از قلم انداخته است.

واضح است که گرچه ارتباط بین امید به زندگی و GDP، مثبت است؛ این ارتباط کاملاً غیرخطی است.

<sup>۱۴</sup> Gross Domestic Product



شکل ۱۱. برازش نمودار به وسیله روش برآوردگر نادارایا - واتسون.

$$w\left(\frac{x_i - x_0}{h}\right) =$$

$$\begin{cases} \left(1 - \left(\frac{|x_i - x_0|}{h}\right)^3\right)^3 & \text{if } \left(\frac{|x_i - x_0|}{h}\right)^3 \leq 1 \\ 0 & \text{if } \left(\frac{|x_i - x_0|}{h}\right)^3 > 1. \end{cases}$$

در شکل ۱۱ (ب) منحنی تابع وزنی مورد استفاده ارائه شده است.

(۳) در گام سوم با به کارگیری تابع وزنی که در گام دوم معرفی کردیم و با بهره‌گیری از داده‌های  $y_i$ ، مقادیر  $\hat{y}_i$  را با استفاده از رابطه‌ی زیر به دست می‌آوریم.

$$\hat{y}_0 = \hat{m}(x_0)$$

(۲) در گام دوم نوبت به انتخاب تابع وزنی می‌رسد. این تابع، باید یک تابع تک مدی متقارن باشد. انتخاب توابع وزنی متفاوت تأثیری در روند کلی کار ندارند و به دلخواه می‌تواند انتخاب شود. در این جا از تابع زیر به عنوان تابع وزنی بهره می‌گیریم. (به عنوان یک تابع وزنی دیگر می‌توانستیم تابع چگالی نرمال را به کار بگیریم.)

تمام مقادیر  $x_i$  ها انجام می دهیم. نمودار شکل ۱۱ (ت) که برآورد تابع رگرسیونی مورد نظر است، حاصل می شود.

$$= \frac{\sum w(x_i, h)y_i}{\sum w(x_i, h)}$$

به شکل ۱۱ (پ) توجه کنید.

(۴) درگام چهارم، مراحل سه گانه ی مذکور را برای

## مراجع

- [۱] رضوی پاریز، ا. (۱۳۸۲)، مقدمه ای بر تحلیل رگرسیون خطی، ترجمه، انتشارات دانشگاه شهید باهنر کرمان.
- [۲] طالبی، ه. (۱۳۷۷)، رگرسیون کلاسیک و مدرن با کاربرد آن، ترجمه، انتشارات دانشگاه شهرکرد.
- [3] Fox, J. (2000). *Nonparametric Simple Regression: Smoothing Scatterplots*, Sage, Thousand Oaks CA.
- [4] Hastie, T.J. and Tibshirani, R.J. (1990). *Generalized Additive Models*, Chapman and Hall, London.
- [5] Linton, O. and Nielsen, J.P. (1995). A kernel method of estimating nonparametric regression based on marginal integration, *Biometrika*, 82, 93-100.
- [6] Nielsen, J.P. and Sperlich, S. (2005). Smooth backfitting in practice. *J. Roy. Statist. Soc.- Series B*, 67, 43- 61.
- [7] Simonoff, J.S. (1996). *Smoothing Methods in Statistics*, Springer, New York.
- [8] Takezawa, K. (2006). *Introduction to Nonparametric Regression*, Wiley, New York.