

## کاربرد رگرسیون چندک در شناسایی شکل توزیع رفاه مورد انتظار جوانان

محمد بامنی مقدم\*، علیرضا خوشگویان فرد\*\*

هدف اصلی این مقاله، آشنا کردن خواننده با کاربرد رگرسیون چندک در تحلیل داده هاست. رگرسیون چندک، رابطه چندک دلخواهی از توزیع متغیر وابسته را با متغیرهای تشریحی از طریق مدل آماری تبیین می کند. در این مقاله، مدل رگرسیون چندک معرفی و به شیوه برآورده پارامترها اشاره می شود؛ به قابلیت شناسایی شکل توزیع که مدل رگرسیون معمولی (میانگین شرطی) آن را دارا نیست، تأکید می شود؛ در پایان با یک مثال عددی از داده های رفاه که در آن براساس یک نمونه تصادفی 684 نفر از جوانان 18 تا 29 سال تهرانی انتخاب شده است، تلاش شده است کاربرد رگرسیون چندک برای رفاه تشریح شود و رابطه رفاه مطلوب جوانان (متغیر وابسته مدل) با تعداد سال های تحصیل آنان (متغیر تشریحی مدل) تحت بررسی قرار گیرد.

---

\* دکترای آمار، عضو هیأت علمی دانشگاه علامه طباطبائی  
<bamenimoghadam@atu.ac.ir>

\*\* فوق لیسانس آمار اقتصادی اجتماعی، مدیر هماهنگی پژوهش های مرکز تحقیقات صداوسیما

کلید واژه‌ها: برآوردهای حداقل مربوعات، حداقل قدر مطلق انحرافات، رگرسیون چندک، رفاه، معیار آنکوک، میانگین شرطی

تاریخ دریافت مقاله: 83/7/27  
تاریخ پذیرش مقاله: 83/10/28

### مقدمه

ابعاد گوناگون پدیده‌ها در آمار، خود را قالب متغیرهای تصادفي نشان می‌دهند و مطالعه آن‌ها با تعیین توزیع (Distribution) آن‌ها میسر می‌شود. برای مثال، وزن نوزادان به عنوان یک پدیده می‌تواند متغیری تصادفي محسوب شود و در صورت مشخص بودن توزیع آن، تعیین این که نوزادی دارای وزن طبیعی است، امکان پذیر می‌شود. توزیع متغیر تصادفي تمام اطلاعات لازم در خصوص آن را به دست می‌دهد، به طوری که می‌توان پدیده‌ها را تفسیر یا پیش‌بینی کرد. به عنوان مثال، با داشتن توزیع وزن نوزادان قادریم پیش‌بینی کنیم که تاچه اندازه ممکن است کودکانی کم وزن داشته باشیم.

در این راستا، معیارهای آماری گوناگونی وجود دارد که هر یک از توزیع متغیر تصادفي اطلاع متفاوتی در اختیار می‌گذارد. برای مثال، واریانس، از نحود پراکندگی و نما، از قلة آن اطلاعی فراهم می‌کند. برخی از معیارها مانند

برجستگی (Kurtosis) کشیدگی یا پхи) و چولگی (Skewness) : کجی) نیز به شکل توزیع اختصاص دارند . شکل توزیع در متغیرهایی مانند درآمد یا هوش از اهمیت ویژه ای برخوردار است . مثلاً اگر توزیع درآمد در جامعه ای چولگی زیادی به چپ داشته باشد، حکایت از وجود افراد محرومی با درآمد بسیار کم دارد . همچنین اگر توزیع نمرات هوش در دانشگاهی متقارن با دم هایی کوتاه باشد، نشان می‌دهد که در مجموع، هوش دانشجویان آن دانشگاه معمولی است و افراد با هوش و کم هوش در آن دانشگاه نسبتاً برابرند.

با توجه به آن‌چه گذشت دور از انتظار نیست اگر ادعا کنیم که استنباط های آماری (Statistical Inference) همگی اطلاعی از توزیع یک متغیر تصادفی به دست می‌دهند. برای مثال، آزمون F در تحلیل واریانس، مقایسه‌ای بین میانگین چند توزیع است؛ محاسبه ضریب همبستگی پیرسون (Pearson) و آزمون مربوط به آن اطلاعی از یک توزیع دو متغیره را فراهم می‌کند؛ مدل‌های رگرسیونی که در ادامه به آن می‌پردازیم نیز مانند دیگر روش‌ها برای بررسی خصوصیات خاصی از توزیع یک متغیر تصادفی به کار می‌روند.

مدل رگرسیون معمولی به تحلیل گر کمک می‌کند تا رابطه میانگین توزیع متغیر تصادفی Y را با

تعدادی متغیر تشریحی بررسی کند. برای روشن شدن مطلب، یک مدل رگرسیون خطی را تنها با یک متغیر تشریحی به این صورت در نظر می‌گیریم:

$$Y_i = \alpha + \beta x_i + \varepsilon_i \quad \text{مدل 1}$$

در این مدل رگرسیونی  $\varepsilon_i$ ها، متغیرهای تصادفی،  $\alpha$  و  $\beta$ ، پارامترهای نامعلوم که باید برآورد شوند و سرانجام  $x_i$ ها مقادیر معلومی از متغیر تشریحی هستند. در صورتی که  $E(\varepsilon_i) = 0$  باشد، آنگاه می‌توان مدل شماره 1 را به صورت دوم بازنویسی کرد:

$$E(Y_i) = \alpha + \beta x_i \quad \text{مدل 2:}$$

کمیت  $E(Y_i)$  را میانگین شرطی (Conditional Mean) متغیر تصادفی  $Y$  مینامند و به همین دلیل آن را با  $E(Y/x_i)$  نیز نشان می‌دهند. بنابراین، مدل شماره 2 بیان می‌کند که میانگین های توزیع  $Y$  در سطوح مختلف متغیر تشریحی در امتداد یک خط راست قرار دارند. به عبارت دیگر، متغیر تصادفی  $Y$  در هر سطح از متغیر تشریحی دارای توزیعی است که میانگین های این توزیع‌ها روی یک خط راست جای گرفته اند. یکی از حالاتی که در آن چنین رابطه‌ای قطعاً برقرار می‌شود زمانی است که دو متغیر  $X$  و  $Y$  دارای توزیع نرمال دو متغیره باشند (باز هم صحبت از توزیع

است!).

از آن جا که میانگین، یکی از معیارهای تمرکز است، آگاهی از آن به تنها ی نمی تواند اطلاعات کاملی از شکل توزیع به همراه داشته باشد. با توجه به این واقعیت، رگرسیون معمولی نیز ممکن است نتواند اطلاعات کافی درباره شکل توزیع متغیر تصادفی تحت مطالعه را - در سطوح مختلف متغیر تشریحی - به دست دهد. چندکها (Quantiles) معیارهای دیگری برای توزیع هستند که «در کنار هم» میتوانند شکل توزیع را جامعتر به تصویر بکشند. برای مثال، اگر دهکهای توزیعی تقریباً دارای فاصله برابری از یکدیگر باشند، انتظار داریم توزیع «نسبتاً» هموار یا یکنواختی داشته باشیم. همچنین اگر دهکهای بالایی دارای فاصله زیاد و دهکهای پایینی دارای فاصله کمی از یکدیگر بشنند، توزیع به سمت راست چوله خواهد بود. اکنون اگر مانند رگرسیون معمولی که برای میانگین به کار می‌رود، یک شیوه رگرسیونی برای چندکها وجود داشته باشد، قادر خواهیم بود شکل توزیع را در سطوح مختلف متغیرهای تشریحی به دست آوریم. این همان هدفی است که رگرسیون چندک دنبال می‌کند.

## 1) معرفی رگرسیون چندک

همان‌طور که در بخش قبل اشاره شد، مدل رگرسیون معمولی، مانند مدل شماره 2، برای میانگین شرطی

برازش داده می‌شود. مدل رگرسیون چندک با ایده ای مشابه برای چندک‌های شرطی (Conditional Quantile) به کار می‌رود. مانند رگرسیون معمولی (میانگین)، کاربردهایی نظیر بررسی رابطه متغیرهای تشریحی با چندک‌ها و هم‌چنین پیش‌بینی آن‌ها برای این نوع از رگرسیون نیز امکان پذیر است. با وجود این، شاید مهم‌ترین کاربرد رگرسیون چندک شناسایی شکل توزیع متغیر وابسته مدل در سطوح گوناگون متغیرهای تشریحی باشد؛ این کار با برآذش مدل‌های رگرسیونی متعدد، به ازای چندک‌های مختلف بر یک مجموعه داده، صورت می‌گیرد.

برای ارائه تعریف دقیقی از مدل رگرسیون چندک ( $\theta \in (0,1)$  ام، ابتدا حالت ساده آن را در نظر  $iid$  می‌گیریم. مدل شماره 1 را با شرط  $(\cdot) \sim F_{x_i}^{iid}$  (تابع  $F$  به یک توزیع دلخواه اشاره دارد) در نظر بگیرید. هدف ما یافتن مدلی است که مثلاً رابطه چندک اول (و نه میانگین) توزیع  $Y$  را با متغیر  $X$  نشان دهد. در این صورت، مدل برای چندک ( $\theta \in (0,1)$  ام متغیر  $Y$  که با  $Q_\theta(Y|x_i)$  نشان داده می‌شود، عبارت است از:

$$Q_\theta(Y|x_i) = \alpha + \beta x_i + F^{-1}(\theta) \quad \text{مدل 3:}$$

تابع فوق، به ازای  $\theta \in (0,1)$ ‌های مختلف، دسته‌ای از خطوط موازی را به دست خواهد داد که دارای عرض

از مبدأهای متفاوتی هستند . در صورتی که  $F(\cdot)$  همان توزیع نرمال (یا هر توزیع متقارن دیگری) باشد، به ازای  $\theta=0.5$  ، مدل شماره 3 همان مدل شماره 2 خواهد بود ، زیرا  $F^{-1}(0.5)=0$  ممکن است به یک تغییر مکان نیاز باشد . اکنون به تعریف کلی مدل رگرسیون چندک می پردازیم . برای این منظور ، فرض کنید  $Y_i = x'_i \beta_\theta + \varepsilon_{\theta i}$  و

$$Q_\theta(Y | x'_i) = x'_i \beta_\theta \quad i=1,\dots,n \quad \text{مدل 4:}$$

که در آن  $x'_i = (1, x_{i1}, \dots, x_{ik})$  و  $\beta'_\theta = (\beta_0, \beta_1, \dots, \beta_k)$  به ترتیب برداری از مقادیر معلوم و پارامترهای نامعلوم بوده و  $\varepsilon_{\theta i}$  یک متغیر تصادفی مشاهده نشدنی است . همچنین ،  $Q_\theta(Y | x_i)$  نمایانگر چندک شرطی  $(0,1) \rightarrow \mathbb{R}$  توزیع  $Y$  است، بنابراین  $Q_\theta(\varepsilon_{\theta i} | x_i) = 0$  . مدل شماره 4 را با شرایط گفته شده ، مدل رگرسیون خطی چندک  $\theta$  می نامیم .

شیوه برآورد پارامترهای مدل رگرسیون معمولی بر حداقل کردن مربع باقیمانده های (انحرافات) مدل مبتنی است که روش حداقل مربعات (Least Squares) نامیده می شود . در این روش، منحنی رگرسیونی به گونه ای برآراش داده می شود که در مجموع، فاصله نقاط از آن به حداقل برسد . در رگرسیون چندک برخلاف رگرسیون معمولی از حداقل کردن مجموع قدر مطلق موزون برای برآورد پارامترهای مدل استفاده

می‌شود که به آن روش حداقل قدر مطلق انحرافات (Least Absolute Deviations) گفته می‌شود. گفتنی است که استفاده از روش LAD که در مدل رگرسیون چندک به کار می‌رود، دارای پیشینه ای طولانی است. در میانه قرن هجدهم، بوسکویچ (Boscovich) یک مدل خطی دو متغیره را بر ای بررسی بیضوی بودن کره زمین از طریق کمینه کردن قدر مطلق خطاهای کار برد. به دنبال آن، لپلاس (Laplace) برآورد ضریب زاویه مدل رگرسیونی بوسکویچ را به طور دقیق معرفی و توزیع مجانبی آن را به دست آورد. ظاهرآ، اج ورث (F. Y. Edgeworth) اولین کسی است که مدل رگرسیونی میانه با چند متغیر تشریحی را در حالت کلی بررسی کرد. توسعه رگرسیون میانه برای هر چندک (Koenker and Bassett) در 1978 صورت گرفت. هدف آن‌ها برآورد بردار پارامترهای  $(\beta_0, \beta_1, \dots, \beta_k)$  در مدل شماره 4 بود که برای این منظور تابع زیانی که در پی می‌آید (قدر مطلق باقیماندهای ای انحرافات موزون) نسبت به عناصر  $\beta_\theta$  کمینه می‌شود:

$$\text{مدل 5: } \varphi_\theta(\beta_\theta) = \sum_i w(\theta) |y_i - x_i' \beta_\theta|$$

$$w(\theta) = \begin{cases} 1 & \theta \leq X_i' \beta_\theta \\ 0 & \text{در این تابع} X_i' \beta_\theta > 1 \end{cases}$$

موزون کردن قدر مطلق باقیماندهای در تابع فوق باعث می‌شود تا خط برآزشی به گونه‌ای باشد که

۱۰۰% داده‌ها تقریباً زیر آن و باقی آن‌ها بالای خط قرار گیرند. کمینه کردن رابطه فوق و یافتن برآورد LAD پارامترها با استفاده از روش‌های برنامه‌ریزی خطی و از طریق بسته‌های نرم افزاری صورت می‌گیرد. در ادامه به ویژگی‌های برآورد LAD اشاره می‌شود.

### ۱-۱) ویژگی‌های حداقل قدر مطلق انحرافات (LAD)

الف) در حالت خاص که مدل تنها شامل عرض از مبدأ و  $\theta=0.5$  است، کمینه کردن رابطه ۵ منجر به کمینه کردن عبارت  $\sum|y_i - \beta_0|$  می‌شود که در این صورت برآورد  $\beta_0$  همان میانه داده‌ها خواهد بود.

ب) برخلاف روش حداقل مربعات، روش حداقل قدر مطلق انحرافات نسبت به داده‌های دور افتاده (Outliers) استوار (Robust) است. این ویژگی ناشی از آن است که برخلاف اهمیت اندازه باقی مانده‌ها در روش حداقل مربعات، در این روش تنها به علامت باقیمانده‌ها توجه نمی‌شود. بنابراین نه تعداد باقیمانده‌هایی که بیشتر (مثبت) یا کمتر (منفی) از چندک مورد نظرند و نه مقدار بزرگی آن‌ها در برآوردها اثرگذار است. پس، داده‌های دور افتاده که تأثیر خود را از طریق بزرگی باقی مانده‌ها نشان می‌دهند، نمی‌توانند برآوردهای LAD را متأثر سازند.

ج) شکل بسته ای برای برآورد پارامترهای این مدل وجود ندارد و از روش‌های عددی برای برآورد آن‌ها استفاده می‌شود. هم‌چنین، جواب‌های نهایی مدل رگرسیون چندک می‌تواند یکتا نباشد. البته یافتن جواب یکتا با انتخاب یک معیار مناسب امکان‌پذیر است. برای مثال، مسئله یافتن میانه 10 عدد را به یاد بیاورید که برای این منظور، میانگین عدد پنجم و ششم به عنوان میانه در نظر گرفته می‌شود. این در حالی است که کلیه اعداد بین عدد پنجم و ششم می‌توانند به عنوان میانه انتخاب شوند. در واقع، با یک قرارداد (معیار)، میانه این اعداد به صورت منحصر به فردی تعیین می‌شود.

د) وقتی  $\varepsilon_{\theta_i}$ ‌ها متغیرهای تصادفی iid باشند، خطوط رگرسیونی به ازای چندک‌های مختلف، موازی خواهند بود (مدل شماره 3 را به یاد بیاورید).

در رگرسیون چندک نیز مانند رگرسیون معمولی (میانگین شرطی) می‌توان استنباط کرد. پاول و بوچنسکی نشان داده‌اند که برآورد LAD پارامترها، سازگار و به طور مجانبی نرمال است (Powell, 1989,. 1998). موضوع اخیر را برای حالت خاص که  $\varepsilon_{\theta_i}$ ‌ها همتوزیع هستند در قضیه زیر بیان می‌کنیم.

## 1-2) قضیه

فرض کنید  $X$  یک ماتریس  $n \times k$  باشد که سطر  $i$  ام آن را  $x'_i$  تشکیل می‌دهد. هم‌چنین تابع چگالی و توزیع

$\epsilon_{\theta_i}$  ها به ترتیب  $f$  و  $F$  باشد. اگر  $V = \lim_{n \rightarrow \infty} [(XX')/n] = F^{-1}(\theta) + \beta_0, \beta_1, \dots, \beta_k$ '،  $g^2(\theta, F) = \frac{\theta(1-\theta)}{[f(F^{-1}(\theta))]^2}$ ،  $f[F^{-1}(\theta)] > 0$  و برآورد آن را با  $b_\theta^*$  نشان دهیم، آنگاه

$$\sqrt{n}(b_\theta^* - \beta_\theta^*) \xrightarrow{d} N[0, g^2(\theta, F)V] \quad \text{مدل 6}$$

اکنون با استفاده از این قضیه می‌توان درباره ضرایب مدل رگرسیون چندک استنباط کرد. بر این اساس، آماره آزمون فرضیه صفر  $H_0: A\beta_\theta = l$  که در آن  $A$  یک ماتریس معلوم با رتبه کامل سطري و  $l$  یک بردار از مقادیر معلوم است، عبارت است از:

$$\lambda = (Ab_\theta - l)' [A(XX')^{-1}A'] (Ab_\theta - l) \hat{g}^{-2} \quad \text{مدل 7}$$

این بخش را با بررسی نحوه شناسایی شکل توزیع و با استفاده از رگرسیون چندک پایان می‌دهیم. برای این منظور ابتدا مدل‌های متعددی به ازای چندک‌های مختلف برداده ها برآرازش داده می‌شود. سپس، برای سطوحی از متغیرهای تشریحی که مذکور است، چندک‌ها براساس مدل‌های برآرازشی پیش‌بینی می‌شوند و از مقایسه چندک‌های پیش‌بینی شده، به شکل توزیع در آن سطوح از متغیرهای تشریحی پی‌می‌بریم. برای مثال، فرض کنید نه دهک از توزیع متغیر وابسته در سطحی از متغیر تشریحی به صورت حالت اول، دوم یا سوم جدول شماره 1 پیش‌بینی شده است.

جدول 1

دهک	0/1	0/2	0/3	0/4	0/5	0/6	0/7	0/8	0/9
حالت اول	3	4	6	7	8	13	19	29	37
حالت دوم	3	8	11	15	19	23	27	32	37
حالت سوم	3	15	22	27	29	31	32	34	37

الف) حالت اول به توزیع چوله به راستی اشاره دارد. به فاصله بیشتر میان دهک‌های بالایی در مقایسه با دهک‌های پایینی توجه کنید. این نشان می‌دهد که 40 درصد از مقادیر بزرگ‌تر متغیر در فاصله وسیعی جای گرفته اند در حالی که 60 درصد باقی‌مانده از مقادیر کوچک‌تر، در فاصله کوتاه‌تری قرار دارند.

ب) حالت دوم نمایان‌گریک توزیع نسبتاً هموار است، زیرا فاصله دهک‌ها تقریباً برابر است  
ج) حالت سوم وضعیتی عکس حالت اول دارد؛ یعنی توزیع چوله به چپی را نشان می‌هد.

آنچه از طریق پیش‌بینی چندک‌ها با کمک مدل به دست آمد - در حالتی که تنها یک متغیر تشريحی در مدل وجود داشته باشد، به کمک رسم خطوط مدل‌های برآزشی نیز دست یافتنی است. در واقع، بررسی

فاصله خطوط رگرسیونی چندک‌های مختلف نیز می‌تواند فاصله چندک‌ها را از یکدیگر در سطوح مختلف متغیر تشریحی نشان دهد.

## (2) کاربرد عملی

این بخش به ارائه یک مثال عددی از رگرسیون چندک اختصاص دارد. یازده مدل رگرسیون چندک همراه با رگرسیون معمولی برای بررسی رابطه رفاه درخواستی افراد (متغیر وابسته مدل) با تعداد سال‌های تحصیل آنان (متغیر تشریحی مدل) به کار می‌رود. رفاه مطلوب و تعداد سال‌های تحصیل به ترتیب با INDEX و EDU نشان داده خواهند شد. داده‌ها به یک نمونه تصادفی 684 نفری از جوانان 18 تا 29 سال تهرانی اختصاص دارد که با روش نمونه‌گیری خوش‌ای سه مرحله‌ای (Three-stage Cluster Sampling) در سال 1381 گردآوری شده است. شایان ذکر است که رفاه درخواست شد. با شاخصی که حاصل از ترکیب 33 سؤال یک پرسشنامه است، سنجیده شده است. گفتنی است مقدار این شاخص از 0 تا 4 تغییر می‌کند، به طوری که هرچه مقدار آن بیشتر می‌شود بر انتظارات بیشتری نیز دلالت دارد.

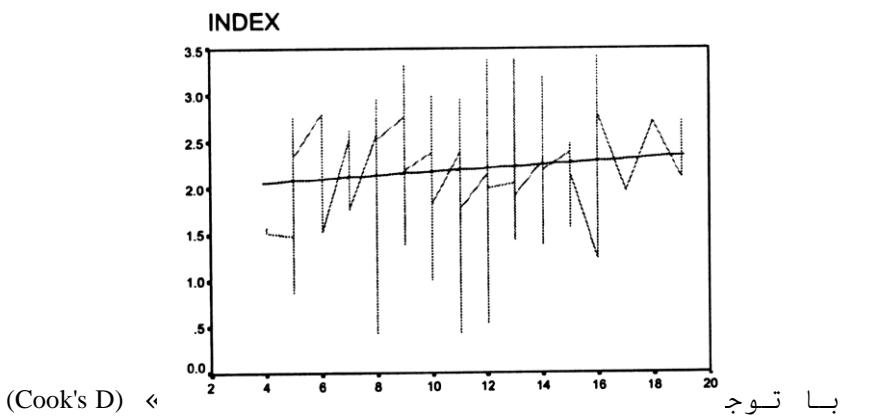
کار را با برآذش یک مدل رگرسیون خطی معمولی با روش حداقل مربعات بر داده‌ها آغاز می‌کنیم. برآورد پارامترهای مدل (عرض از مبدأ و ضریب متغیر تعداد سال‌های تحصیل) در جدول شماره 2

دیده می‌شود؛ بر این اساس مدل برازشی عبارت است از:

$$E(\hat{INDEX}_i) = 1.99 + 0.0189EDU_i$$

که در آن  $E(\hat{INDEX}_i)$  برآورد میانگین توزیع رفاه درخواستی به ازای  $EDU_i$  سال تحصیل است. نمودار پراکنش (Scatter Plot) داده‌ها همراه با خط برازش داده شده در شکل شماره 1 ارائه شده است.

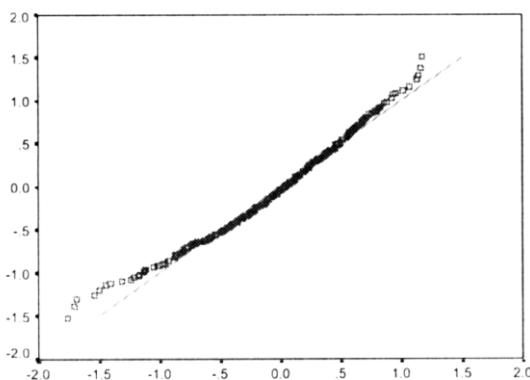
شکل 1: خط رگرسیونی برازشی و نقاط مشاهده شده



به وجود بعداد زیادی داده دورافتاده اشاره داشت. همچنین ترسیم نمودار (Quantile-Quantile Plot) Q-Q در شکل شماره 2 برای بررسی نرمال بودن توزیع باقیمانده‌های مدل، انحراف از توزیع نرمال را نشان می‌دهد. آزمون کلموگرف – اسمیرنوف (Kolmogorov Smirnov Test) نیز فرضیه نرمال بودن این توزیع را رد

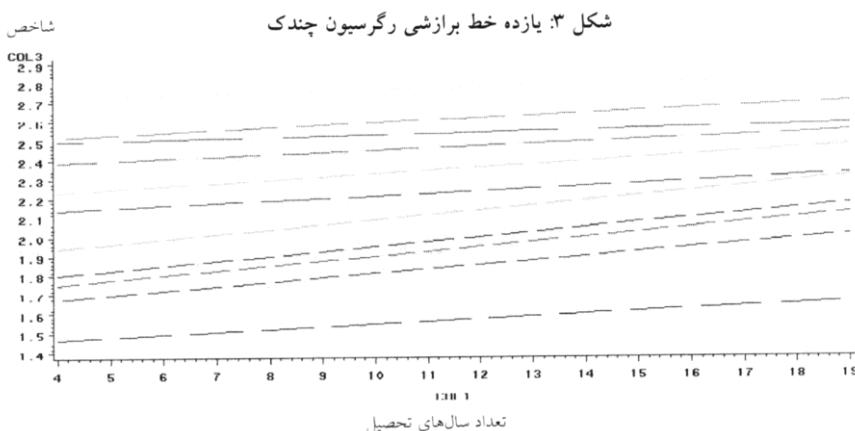
می‌کند. تمام این موارد حکایت از نامناسب بودن مدل رگرسیون معمولی دارد.

شکل 2: نمودار Q-Q برای باقیمانده‌های مدل رگرسیونی  
برآزشی



اکنون به برآزش مدل‌های رگرسیون چندک به ازای یازده چندک مختلف به کمک برنامه ای که در محیط IML از نرم افزار SAS نوشته شده است، می‌پردازیم (برنامه در پیوست ارائه شده است). این برنامه برای مدل 3 و براساس الگوریتمی است که باست و کونکر در 1982 تهیه کردہ‌اند (Bassett & Koenker, 1982). گفتنی است وقتی این چندک، همان میانه باشد، الگوریتم سریع‌تری از طرف مدن و نیلسن پیشنهاد شده است که در SAS/IML در قالب روال LAV (Routine)

استفاده می‌شود (Madsen & Nielsen, 1993). در مجموع، یازده مدل به ازای یازده چندک 0/1، 0/2، 0/25، 0/3، 0/4، 0/5، 0/6، 0/7، 0/75، 0/8، 0/9 بر داده‌ها برآورده شد. برآورد پارامترهای این مدل‌ها در جدول شماره 1 ارائه شده است. شکل شماره 3 نیز خطوط برآورده شده را نشان می‌دهد (خطوط از پایین به بالا مربوط به چندک‌های 0/1 تا 0/9 هستند).



بنابراین، مثلاً، مدل برآورده شده برای چندک 0/2 عبارت است از

$$\hat{Q}_{i,0.2} = 1.58 + 0.0225EDU_i$$

که در آن  $\hat{Q}_{i,0.2}$  برآورده چندک 0/2 توزیع رفاه

درخواستی، به ازای  $EDU_i$  سال تحصیل است (در شکل شماره ۳، خط دوم از پایین، مربوط به این مدل است). پس برآورد چندک  $0/2$  برای افراد با ۱۹ سال تحصیل برابر است با

$$2.0075 = 1.58 + 0.0225 \times 19$$

بنابراین میتوان انتظار داشت که ۲۰ درصد از افراد با ۱۹ سال تحصیل، دارای شاخص انتظارات کمتر از  $2/0075$  و ۸۰ درصد بیش از آن باشند. به همین ترتیب براساس مدل چندک  $0/9$ ، انتظار داریم ۹۰ درصد از افراد با ۱۹ سال تحصیل دارای شاخص انتظارات کمتر از  $2/8249$  و ۱۰ درصد بیش از آن باشند. پس تقریباً ۷۰ درصد از این افراد دارای شاخصی بین  $2/0075$  و  $2/8249$  هستند. توجه کنید که چنین تحلیل‌هایی تنها با مدل‌های رگرسیون چندک قابل انجام است و مدل‌های رگرسیون معمولی چنین قابلیت‌هایی را ندارند. اکنون به یافته‌های حاصل از برآذش این مدل‌ها میپردازیم:

## 2-1) یافته‌های حاصل از برآذش

الف) شکل شماره ۱ حاکی از آن است که خط رگرسیون نمی‌تواند پیش‌بینی‌کننده مناسبی برای شاخص انتظارات باشد. زیرا با توجه به پراکنندگی زیاد داده‌ها در برخی از سطوح تعداد سال‌های

تحصیل، میانگین نمیتواند این شاخص را برای این سطوح به خوبی پیش بینی کند. برای مثال در این شکل افرادی را با 8 تا 12 سال تحصیل ملاحظه کنید. ب) مثبت بودن شبیه خطوط در شکل شماره 3 یعنی ضریب تعداد سال های تحصیل، نشان دهنده رابطه مستقیم بین متغیر وابسته و تشریحی است. بنابراین، با افزایش تعداد سال های تحصیل، مقدار هر یک از یازده چندک شاخص انتظارات نیز افزایش مییابد. بر اساس چندک های برازشی، سال های تحصیل بر چندک های پایینی بیشتر از چندک های بالایی اثر دارد.

ج) برای افراد با تحصیلات بیش تر، فاصله کم چندک های بالایی (چندک 0/9، 0/8، 0/75 و 0/6) در مقایسه با چندک های پایینی نشان میدهد که فشردگی داده ها در بخش بالایی زیاد است. به عبارت دیگر، در سمت راست توزیع شرطی، فشردگی بیش تری در مقایسه با سمت چپ وجود دارد. بنابراین، با افزایش تعداد سال های تحصیل، توزیع شرطی شاخص مطلوب، چوله به چپ می شود.

د) میتوان پیش بینی کرد که مثلاً 70 درصد (از دهک 0/2 تا 0/9) افراد با 19 سال تحصیل دارای شاخصی بین 2/8249 و 2/0075 هستند، در حالی که این فاصله برای افرادی با 4 سال تحصیل از 1/67 تا 2/7184 است.

**جدول 1: مشخصات مدل‌های رگرسیونی**

مجموع قدرمطلق موزون خطاه	برآورد پارامترهای مدل			مدل
	تعداد سالهای تحصیل	عرفه از مبدأ		
—	0/0189	1/99		رگرسیون معمولی
66	0/0129	1/42	0/1	
99	0/0225	1/58	0/2	
110	0/025	1/65	0/25	
118	0/025	1/7	0/3	
127	0/025	1/87	0/4	
128	0/0129	2/09	0/5	
122	0/0166	2/16	0/6	
108	0/0113	2/34	0/7	
98	0/0063	2/47	0/75	
86	0/0125	2/47	0/8	
52	0/0071	2/69	0/9	

### (3) نتیجه‌گیری

این مقاله نشان داد که رگرسیون چندک نه تنها می‌تواند جانشین مناسبی برای رگرسیون میانگین باشد (با جانشین کردن میانه به جای میانگین)، بلکه در برخی از حالات، اطلاعات بیشتری (شكل توزیع) را در مقایسه با رگرسیون میانگین در

اختیار تحلیلگر قرار میدهد. در بخش قبل دیده شد که رگرسیون میانگین به سبب وجود داده های دورافتاده و انحراف از متعادل بودن و هم چنین پراکندگی زیاد متغیر پاسخ در برخی از سطوح متغیر تشریحی، از اعتبار لازم برخوردار نبود؛ در حالی که رگرسیون چندک، یافته های مفیدی را به دست داد.

1. Bassett, G. and Koenker, R. (1982). "**An Empirical Quantile Function for Linear Models with iid Errors**". Journal of American Statistical Association, Vol. 77, No.378: 407-415.
2. Buchinsky, M. (1998). "**Recent Advances in Quantil Regression Models: A Practical Guideline for Empirical Research**". Journal of Human Resources, 33(1): 88-126.
3. Koenker, R. and Bessett, G. (1978). "**Regression Quantiles**". Econometrica, 46: 33-50.
4. Madsen, K. And Nielsen, H. B. (1993). "**A Finite Smoothing Algorithm for Linear L1 Estimation**". SIAM Journal Optimization, Vol. 3: 223-235.
5. Powell, J. (1989). "**Least Absolute Deviation Estimation of the Censored Regression Model**", Journal of Econometrics, 25: 303-325.