

روش های مناسب برای شناسایی داده های پرت و نحوه ی برخورد با آن در آمارگیری های از هزینه و  
درآمد خانوار

مجموعه ی طرح: مجتبی کجفلی

سمانه افتخاری

شهاب جولانی

احسان بهرامی سامانی

گروه پژوهشی طرح های فنی در روش های آمار

پژوهشگده ی آمار



## پیشگفتار

تحلیل رگرسیونی یک ابزار مهم آماری است که به طور معمول در اکثر علوم مورد استفاده قرار می‌گیرد. اغلب، فرض‌های مورد علاقه‌ای بر روی ساختار رگرسیونی به منظور تسهیل ساختن و محاسبه مدل در نظر گرفته می‌شود. برخی از فرض‌های معمولی که بر رگرسیون نهاده می‌شود شامل نرمال بودن متغیر پاسخ و همگونی واریانس پاسخ است. به همین دلیل بیشتر اوقات پیش از مدل بندی داده‌ها، نیازمند استفاده از انواع تبدیلات بر روی پاسخ مورد علاقه به منظور رسیدن به فرض نرمال بودن هستیم (باکس و کاکس، ۱۹۶۴). همچنین همگونی واریانس نیز بایستی پیش از ورود به مرحله تحلیل داده‌ها مورد بررسی قرار گیرد. یک روش مفید می‌تواند بررسی تغییرات پاسخ در سطوح مختلف متغیرهای کمکی مدل باشد (کوک و ویسبرگ، ۱۹۸۳).

از میان روش‌های مختلف برازش مدل رگرسیونی، روش کمترین توان‌های دوم ( $LS$ ) به دلیل قدمت و راحتی محاسبات عموماً مورد استفاده قرار می‌گیرد. اما، امروزه آگاهی بیشتری نسبت به خطرات رخداد نقاط دور افتاده در برازش مدل حاصل شده است (راسو و لروی، ۱۹۸۷). نقاط دور افتاده در داده‌های واقعی بسیار رخ می‌دهند، و اغلب از کنار این نقاط با بی توجهی گذر می‌شود چرا که امروزه داده‌ها توسط رایانه‌ها بدون توجه دقیق پردازش

می‌گردد. نه تنها متغیر پاسخ می‌تواند دورافتاده باشد بلکه متغیرهای تبیینی مدل نیز می‌توانند عامل دورافتادگی قلمداد شوند که به آنها نقاط نافذ گفته می‌شود. هر دوی این گونه نقاط دورافتاده و نافذ می‌تواند موجب عدم کارایی یک تحلیل کمترین توان‌های دوم معمولی گردند. اغلب، چنین نقاط پر نفوذی از چشم کاربر آمارپنهان می‌مانند به دلیل اینکه این نقاط از مانده‌های روش  $LS$  قابل تشخیص نیستند.

دو روش مفید به منظور مقابله با نقاط دورافتاده و نافذ وجود دارد، که یکی ابزارهای تشخیصی در رگرسیون (کوک، ۱۹۷۷ و ۱۹۷۹؛ بلسلی و همکاران، ۱۹۸۰) و دیگری رگرسیون استوار می‌باشد (هابر، ۱۹۷۳ و ۱۹۸۱؛ راسو، ۱۹۸۴؛ راسو و یوهای، ۱۹۸۴). این دو روش دارای هدفی یکسان اما عملکردی در خلاف جهت یکدیگر هستند. در ساختار رگرسیون تشخیصی ابتدا یک رگرسیون برازش داده شده و سپس تحقیق برای تشخیص نقاط دورافتاده بالقوه انجام می‌گیرد تا مدل دوباره با استفاده از داده‌های خوب برازش داده شود در حالی که در رگرسیون استوار، مدلی که مناسب اکثریت داده‌ها باشد برازش داده می‌شود و سپس مشاهداتی که دارای مانده‌های بزرگ از این مدل استوار باشند به عنوان نقاط دورافتاده شناسایی می‌گردند.

همچنین می‌توان تحلیل تاثیر موضعی را به عنوان روشی دیگر به منظور بررسی حساسیت برازش مدل نسبت به فرض‌های مختلف استفاده نمود. با استفاده از این روش نیز می‌توان از طریق حذف یا کم وزن نمودن نقاط مشکوک به شناسایی نقاط دورافتاده یا نافذ با رسم نمودار تاثیر پرداخت. در پژوهش حاضر روشی جدید برای استفاده از خمیدگی نرمال جابه‌جایی درست‌نمایی به عنوان ابزار انتخاب و مقایسه میان رگرسیون‌های استوار مختلف معرفی می‌گردد.

در فصل ۱ به مرور سوابق و ضرورت بحث و بررسی نقاط دور افتاده در داده‌های طرح هزینه و درآمد خانوار می‌پردازیم. فصل ۲ مروری بر مدل‌های رگرسیونی، تبدیلات مناسب و مانده‌های متداول به عنوان ابزارهای رگرسیون تشخیصی خواهد داشت. همچنین در این فصل آماره‌های تشخیص نفوذ و تحلیل تأثیر موضعی مورد بررسی قرار می‌گیرد. انواع رگرسیون استوار نیز در فصل ۳ بررسی می‌شود. در فصل ۴، مدل‌های متغیر پاسخ محدود معرفی می‌شوند. در فصل ۵، داده‌های طرح هزینه و درآمد خانوارهای شهری مورد استفاده در این پژوهش معرفی و مدل رگرسیون معمولی برازش داده شده و انواع روش‌های تشخیصی برای این مدل مورد بررسی قرار می‌گیرند. همچنین در این فصل انواع رگرسیون‌های استوار برای این داده‌ها برازش داده می‌شود و با استفاده از تحلیل تأثیر موضعی با یکدیگر مقایسه می‌شوند. در گروه مطالعاتی مذکور، جناب آقای دکتر گنجعلی عضو هیأت علمی دانشگاه شهید بهشتی به عنوان مجری طرح و خانم سمانه افتخاری و آقایان شهاب جولانی و احسان بهرامی سامانی در گروه عضویت داشته‌اند که بدین وسیله از زحمات بی‌دریغ یکایک این افراد تشکر و قدردانی می‌شود.

نشریه‌ی حاضر، حاصل تلاش گروه مطالعاتی فوق‌الذکر است که امید است مورد توجه و استفاده‌ی مسئولین مسئولین و کارشناسان ذیربط قرار گیرد.

**گروه پژوهشی طرح‌های فنی و روش‌های آماری**



# فهرست مندرجات

۱	مقدمات	۱
۱	مقدمه‌ای بر طرح آمارگیری هزینه و درآمد و مرور سوابق	۱.۱
۳	۱.۱.۱ مرور سوابق	۱.۱.۱
۶	۲.۱ نقاط دور افتاده	۲.۱
۸	۱.۲.۱ نمودارهای مانده‌ها	۱.۲.۱
۱۱	۳.۱ مثال‌ها	۳.۱
۲۳	۲ مدل‌های رگرسیونی، داده‌های پرت و شناسایی آن‌ها	۲
۲۳	۱.۲ مقدمه	۱.۲

۲۴	.....	مروری بر مدل رگرسیون خطی	۲.۲
۲۴	.....	مدل رگرسیونی خطی ساده	۱.۲.۲
۲۵	.....	بررسی مفروضات مدل رگرسیونی پیش از انجام تحلیل	۳.۲
۲۵	.....	تبدیلات باکس-کاکس	۱.۳.۲
۲۷	.....	ناهمگونی واریانس	۲.۳.۲
۲۸	.....	تحلیل مانده‌ها	۴.۲
		حذف مشاهدات و اثر نقاط پرنفوذ در برآورد $\hat{\beta}$ و مجموع	۱.۴.۲
	۳۳	.....	توان دوم خطاها
۳۴	.....	آماره‌های تشخیص نقاط نافذ	۲.۴.۲
۳۶	.....	آماره‌ی کوک $D$	۳.۴.۲
		معیار تفاوت برآزش‌ها قبل و بعد از حذف آزمودنی	۴.۴.۲
	۳۸	.....	( ${}^1DFFITs$ )
		آماره تفاوت پارامترها قبل و بعد از حذف آزمودنی	۵.۴.۲
	۳۸	.....	( ${}^2DFBETAS$ )
۳۹	.....	آماره مانده حذف	۶.۴.۲

<sup>۱</sup> Difference of Fits<sup>۲</sup> Difference of Betas



۷.۴.۲	آماره نسبت واریانس‌ها قبل و بعد از حذف آزمودنی
۴۰	( ${}^2COVRATIO$ )
۵.۲	تحلیل تأثیر موضعی برای ارزیابی حساسیت مدل نسبت به داده‌های
۴۱	پرت
۱.۵.۲	تحلیل تأثیر موضعی
۴۲	
۲.۵.۲	خمیدگی قائم
۴۴	
۳.۵.۲	مفهوم خمیدگی
۴۴	
۴.۵.۲	محاسبه خمیدگی قائم در رویه هندسی
۴۶	
۵.۵.۲	تأثیر موضعی در مسئله وزن‌های موردی در رگرسیون ساده
۵۰	خطی
۳	مدل‌های استوار
۵۵	
۱.۳	مقدمه
۵۵	
۲.۳	تعاریف اساسی نظریه‌ی رگرسیون استوار
۵۶	
۱.۲.۳	نقطه‌ی فروریزش
۵۸	
۲.۲.۳	کارایی
۵۸	

<sup>۲</sup> Covariance Ratios

۵۹	.....	خاصیت کراندار بودن نقاط نافذ برآوردگرها	۳.۲.۳
۵۹	.....	مدل‌های استوار	۳.۳
۵۹	.....	روش کمترین قدرمطلق انحرافات ( $LAD$ ) <sup>۴</sup>	۱.۳.۳
۶۱	.....	$M$ رگرسیون	۲.۳.۳
۶۳	.....	$M$ رگرسیون تعمیم یافته	۳.۳.۳
۶۵	.....	رگرسیون کمترین میانه‌ی توان‌های دوم ( $LMS$ )	۴.۳.۳
۶۵	.....	رگرسیون کمترین توان‌های دوم پیراسته ( $LTS$ ) <sup>۵</sup>	۵.۳.۳
۶۶	.....	$S$ برآورد	۶.۳.۳
۶۸	.....	$MM$ برآوردگرها	۷.۳.۳

#### ۴ مدل‌های متغیر پاسخ محدود ۷۰

۷۰	.....	مقدمه	۱.۴
۷۱	.....	مدلهای رگرسیون سانسوریده یا بریده	۲.۴
۷۷	.....	مدل توییت	۱.۲.۴
۸۲	.....	مدل رگرسیون گزینش شده	۲.۲.۴
۸۸	.....	مدل رگرسیون بریده	۳.۲.۴

<sup>۴</sup> Least-Absolute Deviations

<sup>۵</sup> Least Trimmed of Squares

۳.۴	مانده‌ها در مدل‌های متغیر پاسخ محدود	۹۰
۵	کار کاربردی بر روی داده‌های هزینه و درآمد	۹۵
۱.۵	داده‌های طرح هزینه و درآمد سال ۱۳۸۵	۹۵
۲.۵	برازش مدل رگرسیون خطی	۹۹
۱.۲.۵	بررسی تشخیصی مدل رگرسیون خطی	۱۰۱
۳.۵	نتایج رگرسیون استوار بر روی داده‌های طرح هزینه و درآمد خانوار	۱۰۹
۱.۳.۵	تحلیل تاثیر موضعی برای داده‌های طرح هزینه و درآمد	۱۱۴
A	برنامه‌های محاسباتی مربوط به مدل‌بندی داده‌های طرح هزینه	
	و درآمد خانوار شهری	۱۱۶
	مراجع	۱۳۲



## مقدمات

### ۱.۱ مقدمه‌ای بر طرح آمارگیری هزینه و درآمد و مرور سوابق

طرح آمارگیری هزینه و درآمد خانوارهای شهری و روستایی یکی از طرح‌های آماری گسترده در کشور می‌باشد که با قدمتی ۳۸ ساله یکی از قدیمی‌ترین و مهم‌ترین طرح‌های آمارگیری اجرا شده توسط مرکز آمار ایران است. مجریان این طرح هر ساله با مطالعات و بررسی‌های لازم، تطبیق آن را با استانداردها و توصیه‌های بین‌المللی، بررسی و تغییرات مورد نیاز در راستای تکامل آن را تهیه و اجرای آن را به کارشناسان مربوط واگذار می‌کنند. هدف کلی این طرح، برآورد میانگین هزینه‌های خوراک، غیرخوراک و هزینه کل و نیز برآورد میانگین درآمد برای یک خانوار در سطح نقاط شهری و روستایی هریک از استان‌ها و کل کشور می‌باشد. به علت وسعت پوشش مکانی و نیز نوع اقلام مورد پرسش، نتایج حاصل از این طرح بسیاری از اطلاعات مورد نیاز اقتصادی همچون شناسایی الگوی مصرف، توزیع درآمد و نیز میزان و تحولات بر خورداری خانوارها از امکانات و تسهیلات زندگی و اجتماعی را در مقیاس داخلی و

بین‌المللی فراهم می‌آورد که همه اینها در سیاست‌گذاری‌های کلان اقتصادی دولت و بخش خصوصی مفید واقع می‌شود.

در این طرح مشخصات اساسی خانوار از جمله خصوصیات اجتماعی اعضای خانوار که شامل سن، وضعیت سواد (باسواد یا بی‌سواد) و تحصیل (شاغل به تحصیل یا خیر)، وضعیت فعالیت (شاغل، بیکار جویای کار، دارای درآمد بدون کار، محصل، خانه‌دار و سایر) و وضعیت زناشویی (دارای همسر، بی‌همسر بر اثر فوت، بی‌همسر بر اثر طلاق و هرگز ازدواج نکرده) است، ثبت می‌گردد و همچنین مشخصات محل سکونت و تسهیلات و لوازم عمده زندگی مانند نحوه تصرف محل سکونت (ملکی عرصه و اعیان، ملکی اعیان، اجاری، رهن، در برابر خدمت، مجانی و سایر) و نوع سوخت عمده مصرفی خانوار (نفت سفید، گازوئیل، گاز، برق و سایر) مورد بررسی قرار گرفته و در سطحی گسترده به بررسی و ثبت هزینه‌های خانوار در قالب دو بخش خوراکی و غیر خوراکی پرداخته شده است. از جمله اقلام مورد توجه در قسمت هزینه‌های خوراکی می‌توان به هزینه‌های مربوط به نان، گوشت، لبنیات و خشکبار مصرفی خانوار اشاره نمود. در بخش هزینه‌های غیر خوراکی به هزینه‌های پوشاک، مسکن و تسهیلات آن، هزینه‌های بهداشتی و درمانی خانوار، هزینه‌های حمل و نقل خانوار، هزینه‌های ارتباطات، خدمات فرهنگی و تفریحی، آموزش و تحصیل و سرمایه‌گذاری خانوار توجه گردیده است. همچنین درآمدهای خانوار که شامل مواردی همچون درآمدهای حاصل از مشاغل حقوق‌بگیری و آزاد و درآمدهای متفرقه می‌باشد نیز ثبت گشته است.

با توجه به اهمیت اطلاعات درآمد خانوار در مباحث مربوط به توزیع درآمد، آمارهای فقر و شاخص‌های رفاه اقتصادی، ارزیابی‌ها و سیاست‌گذاری‌های اقتصادی با اتکا به این مهم، ضروری و حائز اهمیت فراوان می‌باشد. همچنین آگاهی از هزینه و درآمد خانوارهای

کشور می‌تواند منعکس‌کننده میزان و ترکیب هزینه‌ها و درآمدها و چگونگی استفاده خانوارها از امکانات و تسهیلات زندگی در سطح ملی و منطقه‌ای باشد. باید توجه داشت که استفاده از مدل‌های مناسب آماری برای پیش‌بینی و بررسی ساختاری اطلاعات بدون در نظر گرفتن خصوصیات توزیعی متغیرهای به کاررفته، اغلب برآوردهای گمراه‌کننده‌ای را نتیجه می‌دهد. در این میان شناسایی داده‌های پرت به دلیل تاثیر قابل ملاحظه‌ای که بر برآزش مدل القا می‌نماید از اهمیت ویژه‌ای برخوردار است.

### ۱.۱.۱ مرور سوابق

#### طرح داده کاوی و کاربرد آن در کیفیت داده‌ها

بر اساس گزارش طرح مطالعاتی داده کاوی و کاربرد آن در کیفیت داده‌ها (این طرح از طریق مرکز آمار ایران قابل دسترسی می‌باشد)، اغلب اطلاعات و داده‌های موجود در طرح هزینه و درآمد خانوار توزیع‌های ناشناخته یا پیچیده‌ای دارند که به راحتی نمی‌توان آن توزیع‌ها را شناسایی نمود و مورد استفاده قرار داد. بنابراین برای تحلیل داده‌ها و اطلاعات فوق استفاده از روش‌هایی که نیاز به دانستن توزیع متغیرها ندارد از اهمیت خاصی برخوردار است. خوشه بندی را می‌توان هم به روش پارامتری و هم ناپارامتری انجام داد که در صورت عدم امکان دسترسی به اطلاعات مورد نیاز برای استفاده از روش پارامتری، می‌توان از خوشه بندی ناپارامتری که با توزیع داده‌های موجود سروکار نداشته و اغلب با استفاده از معیارهای تشابه و عدم تشابه انجام می‌گیرد، استفاده نمود. اغلب در پایگاه داده‌های با حجم بالا تعداد متغیرها، حجم داده‌ها یا هر دو خیلی زیاد است که طرح هزینه و درآمد خانوار نمونه بارزی از آن است، برای خوشه بندی داده‌ها در این پایگاه‌ها یک روش بسیار مفید و ارزشمند تحلیل خوشه‌ای و

تحلیل مناسب در خصوص گروه‌بندی متغیرها تحلیل عاملی است. از آنجا که تعداد متغیرها و همچنین تعداد خانوار نمونه در طرح هزینه و درآمد خانوار زیاد می‌باشد، استفاده از روش‌های مرسوم خوشه‌بندی همچون خوشه‌بندی سلسله‌مراتبی یا دو مرحله‌ای، پیشنهاد نمی‌شود. لذا برای انجام تحلیل خوشه‌ای از الگوریتم کلارا که یکی از روش‌های خوشه‌بندی پم (بخش بندی اطراف مدوئید) است و از جمله روش‌های خوشه‌بندی در داده کاوی است استفاده شده است. این روش‌ها در مقایسه با روش‌های سلسله‌مراتبی نیاز به صرف وقت خیلی کمتر و امکانات سخت افزاری معمولی دارند. از آنجا که در روش خوشه‌بندی تعداد خوشه‌ها را نمی‌توان با استفاده از استدلالی علمی از قبل تعیین کرد، لذا از حداقل انتخاب تعداد خوشه شروع به خوشه‌بندی نموده و این کار را تا آنجا که نمودار سایه نمای مناسب به دست آید (ماکسیمم مقدار میانگین سایه نما) ادامه می‌دهند. (برای راهنمایی بیشتر به گزارش طرح مطالعاتی داده کاوی و کاربرد آن در کیفیت داده‌ها در مرکز آمار مراجعه کنید).

در این طرح برای شناسایی و اصلاح مشاهدات دور افتاده از خوشه‌بندی استفاده گردیده که ابتدا سرجمع‌های هر بخش پرسشنامه را محاسبه کرده و سپس به خوشه‌بندی هزینه‌های مربوط به هر بخش پرسشنامه پرداخته شده است، بدین منظور خوشه‌بندی با استفاده از متغیرهایی چون بعد خانوار، فصل آمارگیری و هزینه هر بخش (شامل بخش هزینه‌های خوراک، پوشاک و ...) و نیز هر ۱۳ بخش به طور کلی در سطح مناطق روستایی کل کشور و استان تهران برای تمامی خانوارهای نمونه فصل ۱ سال ۱۳۸۴ انجام شده است. برای این کار با استفاده از الگوریتم خوشه‌بندی کلارا، خانوارها در بخش‌های هزینه خوشه‌بندی شده‌اند که بر اساس ضریب سایه نما که تعلق داده‌ها به خوشه‌ها را نمایش می‌دهد، بزرگترین میانگین ضریب سایه نما که بیانگر بهترین تعداد خوشه‌هاست انتخاب شده است.



سپس داخل خوشه‌ها برای هر بخش داده‌های دورافتاده شناسایی شده‌اند. داده‌هایی که در خارج از فاصله  $1/5$  برابر دامنه میان چارکی (متغیر مورد نظر در آن بخش) قرار دارند را به عنوان مشاهدات دورافتاده بالقوه معرفی شده‌اند. سپس با استفاده از سایر فراداده‌های موجود و مرتبط با خانوار در فایل داده‌های مربوط به خانوار دورافتاده بودن یا نبودن اطلاعات خانوار تعیین شده است.

یکی از مشکلات روش گفته شده در بالا نیازمندی به فراداده‌هاست که اغلب جمع آوری نشده‌اند یا به طور کامل در دسترس نمی‌باشند. همچنین روش ذکر شده تنها به مسئله داده‌های دورافتاده در داده‌های خام توجه می‌نماید که روش تشخیصی آن نیز بر مبنای در نظر گرفتن توزیعی متقارن (استفاده از نمودار جعبه‌ای) برای متغیرهای مورد نظر است که ممکن است در عمل رخ ندهد. از طرف دیگر مسئله داده‌های دورافتاده زمانی اهمیت بیشتری خواهد داشت که مدل‌بندی و استخراج اطلاعات از داده‌های موجود مورد نظر باشد. این امر خود مستلزم استفاده از تعاریفی دقیق‌تر از مفهوم داده‌های دورافتاده می‌باشد که گاه به آنها داده‌های با نفوذ بالا اطلاق می‌گردد و تمامی مفاهیم گفته شده تنها در قالب مدل‌بندی معنی و مفهوم واقعی خود را می‌یابد. لذا در پژوهش حاضر به مسئله داده‌های پرت زمانی که به مرحله استخراج اطلاعات به منظور پیش‌بینی یا استنباط وضع موجود توجه گردد، اهمیت داده می‌شود و مسئله داده‌های پرت و شناسایی آنها در راستای دستیابی به مدلی بهتر مورد توجه قرار می‌گیرد.

## ۲.۱ نقاط دورافتاده

مسئلهٔ نقاط دورافتاده همواره در برازش یک مدل مورد توجه قرار می‌گیرند و معرفی آنها منوط به داشتن معیاری برای میزان دوری از اکثریت داده‌هاست که بر طبق مانده‌های مرتبط با مدل مورد استفاده تعریف می‌گردند. یکی از آشناترین این مدل‌ها، مدل رگرسیونی ساده می‌باشد که شامل برازش مدل‌هایی به یک پاسخ پیوسته است که امید داریم به مقادیر چندین متغیر کمکی وابسته باشد. همچنین در رگرسیون، فرض می‌کنیم که متغیر پاسخ به همراه خطا مشاهده گردیده در حالی که متغیرهای کمکی مقادیر معلومی بدون خطا هستند (که البته در عمل فرض واقع بینانه‌ای نمی‌باشد) و رابطه‌ی میان این دو مجموعه از متغیرها به مجموعه‌ای از پارامترهای نامعلوم که توسط روش کمترین توان‌های دوم برآورد می‌شوند، وابسته است.

در رگرسیون خطی، متغیر پاسخ به صورت خطی به بردار پارامترهای نامعلوم وابسته است. دو کلاس کلی‌تر از مدل‌ها شامل مدل‌هایی است که غیر خطی در پارامتر هستند و مدل‌های تعمیم‌یافته که در آن‌ها، احتیاجی به پیوسته بودن پاسخ نبوده و وابستگی میان میانگین متغیر پاسخ و متغیرهای مستقل توسط تابع ربط برقرار می‌شود. توجه ما در اینجا بیشتر به مدل‌هایی است که به صورت رگرسیون چندگانه باشند یعنی پاسخ به چندین متغیر کمکی وابسته باشد.

برونلی<sup>۱</sup> (۱۹۶۰)، رگرسیون چندگانه را توصیف کرده و به برآورد پارامترها و آزمون فرض پارامترهای مدل پرداخته است، اما از روش‌های نموداری همچون رسم داده‌ها، یا مانده‌های مدل برازش شده برای تشخیص نقاط دورافتاده استفاده نکرده است. از زمان انتشار کتاب برونلی، تفاوت‌های تحسین برانگیزی در تحلیل‌های آماری رخ داده که با وجود نرم افزارهای قوی گرافیکی آماری به راحتی می‌توان خصوصیات مدل‌های مختلف را با رسم نمودارهای

<sup>۱</sup> Brownlee

گوناگون بررسی کرد.

در این بخش، به روش‌های نموداری برای ساختن و بررسی مدل‌های رگرسیونی توجه خواهیم نمود. یکی از کاربردهای این نمودارها که به نمودارهای تشخیصی<sup>۲</sup> معروف هستند، شامل تشخیص پرت بودن یک یا چند مشاهده می‌باشد. داده‌هایی که دور افتاده هستند، به حالت‌های مختلفی می‌توانند رخ دهند. در زیر چندین حالت آورده شده است.

۱ - خط‌هایی که ممکن است در متغیر پاسخ یا متغیرهای کمکی و یا هر دو وجود داشته باشند. که ممکن است اینگونه خط‌ها ناشی از اشتباه در اندازه‌گیری‌ها یا خط‌های ایجاد شده در ورود داده‌ها باشند.

۲ - خط ناشی از در نظر گرفتن شکل صحیح رابطه ریاضی مدل می‌باشد.

۳ - ممکن است به تبدیل مقیاس مانند لگاریتم پاسخ یا دیگر تبدیل‌ها، نیاز باشد (باکس و کاکس<sup>۳</sup>، ۱۹۶۴).

۴ - قسمت سیستماتیک مدل و مقیاس آن ممکن است درست بوده اما توزیع خط‌های پاسخ متفاوت از نرمال باشند، برای مثال دم‌هایی کلفت تر از توزیع نرمال داشته باشد، که در نتیجه استفاده از روش ماکسیمم درست‌نمایی برازش مدل را محدود می‌سازد (دراپر و اسمیت، ۱۹۶۶).

یک منبع بالقوه‌ی خط‌های بزرگ (در مطالعاتی که کمیت داده‌ها کوچک باشد)، حذف ارقام اعشاری داده‌هاست که باید با دقت زیادی انجام گیرد و حتی الامکان از نرم‌افزارهایی استفاده شود که پیش‌زمینه‌ی خود کار حذف اعشار را نداشته باشند، تفاوت ویژه‌ی میان داده‌ها و مدل، تعیین‌کننده‌ی آسان یا سخت بودن تشخیص ناهنجاری در مدل است.

<sup>۲</sup> Diagnostic Plots

<sup>۳</sup> Box and Cox

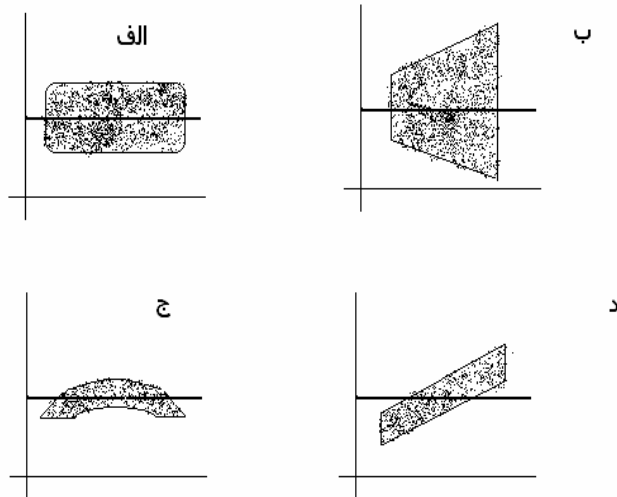
در این قسمت به منظور مطالعه مانده‌ها و نقاط پرنفوذ به بررسی مشاهدات عجیبی می‌پردازیم که بر اساس یکی از سه مورد اول ذکر شده در بالا ناشی شده باشند. بدین منظور در زیر مطالب ساده ولی بسیار مهمی را در ردیابی این قبیل داده‌ها به روش گرافیکی مرور می‌کنیم. مبانی نظری مربوط به این روش‌ها در فصل ۲ گزارش آورده شده است. حالت ۴ نیز در روش‌های رگرسیون استوار در فصل ۳ بررسی خواهد شد.

### ۱.۲.۱ نمودارهای مانده‌ها

نمودارهای بسیاری برای بررسی تفاوت‌های ساختاری مدل با داده‌ها استفاده می‌شود. این نمودارها بر اساس مانده‌ها رسم می‌گردند. در یک مدل به طور کلی، مانده به تفاضل مقدار برازش شده تحت مدل از مقدار مشاهده شده گفته می‌شود. نمودارهایی که در کنترل فرض‌های مدل و بررسی انحرافات از مدل، مناسب می‌باشند به قرار زیرند.

#### ۱ - نمودار دنباله زمانی مانده‌ها

برای بررسی وجود تاثیر گذشت زمان (در این حالت فرض می‌کنیم که زمان به عنوان یک متغیر مستقل در مدل استفاده شود) در مقادیر مانده‌های مدل از نمودار مانده‌ها در مقابل زمان استفاده می‌گردد. اگر از دور به نمودار زمانی نگاه کنیم و تجسمی از یک مرز افقی برای مانده‌ها بدست آوریم که می‌توان آن را به صورت قسمت الف از شکل ۱.۱ نشان داد، این نشان می‌دهد که اثر زمانی در دراز مدت روی داده‌ها موثر نیست. اما اگر شکل مانده‌ها شبیه یکی از اشکال ب یا ج یا د در نمودار ۱.۱ باشد، نتیجه می‌گیریم که اثر زمانی در این حالات، در مدل به حساب نیامده است. اگر رفتاری مانند نمودار ب در شکل ۱.۱ رخ دهد می‌توان نتیجه گرفت که واریانس ثابت نبوده بلکه با زمان افزایش می‌یابد.



شکل ۱.۱: حالات ممکن نمودارهای مانده‌ها

نمودار د در شکل ۱.۱ نشان دهنده‌ی نیاز به یک جمله‌ی خطی برحسب زمان در مدل است و نمودار ج در شکل ۱.۱ گویای این مطلب است که باید جملات خطی و درجه‌ی دوم برحسب زمان به مدل اضافه شوند. البته ترکیب‌ها و تغییراتی از این اثرها نیز می‌توانند رخ دهند (مثلاً در مورد نمودار د شکل ۱.۱ شیب در جهت مخالف باشد و نظایر آن).

## ۲ - نمودار مانده‌ها نسبت به مقادیر برازش یافته:

با رسم مقادیر مانده‌ها در مقابل مقادیر برازش یافته، باز مانند نمودار دنباله‌ی زمانی می‌توان انحرافات مدل را بررسی نمود.

۱ - اگر نمودار همانند نمودار الف شکل ۱.۱ یک مرز افقی را نشان دهد، درست بودن تحلیل‌های کمترین توان‌های دوم (و نرمال بودن در مدل رگرسیون خطی) تایید می‌گردد. در مقابل اگر اشکال ب یا ج یا د از شکل ۱.۱ رخ دهند نتایج به صورت زیر خواهند بود: