

فضل ششم: عیب شناسی در رنگ‌سینون چندگانه و انجام تبدیلات بر روی آن
در این فصل روش‌های را برای آزمون فرضیات اساسی در مدل رنگ‌سینون چندگانه و نهایتاً عیب شناسی
در این مدل ارائه خواهیم نمود. همچنین در صورت امکان روش‌های را برای غلبه بر این مشکلات
(که از سبب ثابت بودن واریانس خطاها و غیره می‌باشد) بیان می‌کنیم و مسئله را
پیشتر خواهیم نمود.

۹.۱ عیب شناسی در رنگ‌سینون چندگانه

در هنگام برآزش یک مدل رنگ‌سینون چندگانه برای بررسی به‌قراری فرضیات اساسی توجه به نکات زیر
ضروری است:

- ۱- مشخص کنیم که آیا مدل برآزش شده به داده‌ها معیار است یا خیر (برآزش کافی به داده‌ها دارد یا خیر)
ابزار اصلی جهت رسیدن به این هدف و بررسی به‌قراری فرضیات رسم نمودار ماندن‌های استاندارد و
مقایسه برآزش شده است. همچنین رسم نمودارهای پراکنش حاشیه‌ای نیز در انجام این امر بسیار مهم است.
- ۲- مشخص نمودن آنکه کدام یک از نقاط دارای تأثیر زیادی بر برآزش مدل است یا نه (معمولاً در این کار می‌تواند به کمک کویچ‌ها انجام گردد).
- ۳- مشخص نمودن داده‌های پرت در صورت وجود. یعنی داده‌هایی که از مخرج کلی مدل برآزش
شده و سایر داده‌ها تبعیت نمی‌کنند.
- ۴- ارزیابی تأثیر حرکت از متغیرهای پیشگویی‌کننده بر روی متغیر پاسخ یا استناد از نمودارهای متغیر
اچانه شده (قسمت ۳.۲.۴).
- ۵- ارزیابی مقدار همبستگی بین متغیرهای پیشگو، استناد از معیارهای معرّفی شده
تعبیرات مثال عوامل تورم واریانس.
- ۶- آزمون آنکه آیا فرضیه ثابت بودن واریانس خطاها به‌قراری است یا خیر. اگر خیر چگونه و در آن
برای این مشکل غلبه نمود.

۷- اگر داده‌های زمانی داریم یعنی داده‌ها در طول زمان جمع‌آوری شده‌اند بررسی اینکه این داده‌ها بر روی زمان همبستگی دارند یا خیر.

یادآوری رگرسیون چندگانه در فرمت ماتریسی: فرض کنید بردار متغیر پاسخ $Y_{n \times 1}$ بردار متغیر پیش‌بینی $X_{n \times (p+1)}$

$$Y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}, \quad X = \begin{pmatrix} 1 & x_{11} & \dots & x_{1p} \\ \vdots & x_{21} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \dots & x_{np} \end{pmatrix}$$

ماتریس طرح باشد، یعنی

بر اساس ماتریس X و بردار Y ، پارامترها و خطای تصدیق زیر قابل تعیینند:

$$\beta = (\beta_0, \beta_1, \dots, \beta_p)' \quad , \quad \epsilon = (\epsilon_1, \epsilon_2, \dots, \epsilon_n)'$$

پس برای β می‌توان نوشت:

$$Y = X\beta + \epsilon \quad , \quad \epsilon \sim (0, \sigma^2 I)$$

$$\hat{\beta} = (X'X)^{-1} X'Y \Rightarrow \hat{Y} = X\hat{\beta} = \underbrace{X(X'X)^{-1} X'}_H Y = HY \quad , \quad H = X(X'X)^{-1} X'$$

$$e = Y - \hat{Y} = Y - HY = (I - H)Y$$

۱.۱.۴. لوریج h_{ii} در رگرسیون چندگانه: همانطور که دیدیم در مدل رگرسیون خطی ساده، لوریج h_{ii} معیاری برای نشان دادن تأثیرات داده نام بردار مدل بر آنش شده بودند. در مدل رگرسیون چندگانه نیز این مقادیر همین نقش را داشته و بصورت زیر تعریف می‌شوند:

$$\hat{y}_i = HY \quad , \quad \hat{y}_i = \left(\frac{\text{مجموع نام ماتریس}}{H} \right) \times Y = \sum_{j=1}^n h_{ij} y_j = h_{ii} y_i + \sum_{j \neq i} h_{ij} y_j$$

پس برای \hat{y}_i می‌توان نوشت: $h_{ij} = x_i (X'X)^{-1} x_j'$ ، $h_{ii} = x_i (X'X)^{-1} x_i'$ ، $h_{ii} > 0$ ، $\frac{p+1}{n} < h_{ii} < 1$ ، h_{ii} نشان دهنده تأثیرات داده نام بردار است.

بنابراین h_{ii} سطر نام ماتریس طرح یعنی $(x_{i1} \ x_{i2} \ \dots \ x_{ip})$ و h_{ii} قطر و نام آن است.

۲.۱.۴. بررسی خواص فائده‌ها در رگرسیون چندگانه: بردار فائده‌ها در مدل رگرسیون چندگانه بصورت زیر تعریف می‌شود:

۱۸
$$e = y - \hat{y} = (I - H)y \Rightarrow E(e) = E(y - \hat{y}) = E[(I - H)y]$$

$$\Rightarrow E(e) = (I - H)E(y) = (I - H)X\beta = X\beta - HX\beta = 0$$

$$Var(e) = Var[(I - H)y] = \sigma^2(I - H)I(I - H)' = \sigma^2(I - H)$$

زیرا H یک ماتریس متقارن و $(I - H)$ خودتوان است. $H = H'$

$$(I - H)(I - H)' = I - 2H + H^2 = I - 2H + H = I - H$$

$$Cov(e, \hat{y}) = \sigma^2(I - H) = 0$$

بنابراین هم‌توان ماندن آنها را استاندارد را بصورت زیر تعریف نمود:

$$Var(e_i) = \sigma^2(1 - h_{ii}) \Rightarrow r_i = \frac{e_i}{\sqrt{1 - h_{ii}}}, \quad s = \sqrt{MSE} = \sqrt{\frac{SSE}{n - (p + 1)}}$$

چنانچه r_i از ۲ تا ۱۳ باشد. در توان گفت که مقدار نام یک داره پرت است و در غیر این صورت

پرت نیست. به توجه داشته که در صورتی یک مقدار پرت خواهد بود که مدل برازش شده معیار باشد و

در غیر این صورت نباید بر حسب داره پرت به آن زود. همچنین ممکن است برای یک مجموعه

(پرتی و غیره)

از داره ها $r_i = -1.8$ شود که در مقایسه با فرآیند داره ها آن داره یک راده پرت به نظر

آید. بنابراین قبل از حذف یک داره بعنوان داره پرت، بهترین کار نگاه به نمودار پراکنش داره ها

و بررسی چند معیار بطور همزمان برای پرت بودن داره هاست که در ادامه معرفی خواهیم کرد.

استفاده از مانده ها و مانده های استاندارد برای بررسی مدل: به زبان ساده، مدل رگرسیون

چندگان یک مدل معتبر برای داره هاست هرگاه میانگین شرطی y ، شرط x یک تابع خطی از

x باشد و واریانس شرطی $Var(y|x=x) = \sigma^2$ ثابت باشد.

$$E(y|x=x) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$$

بنابراین یک مدل رگرسیونی معتبر است که نمودار مانده ها را استاندارد در مقابل ترتیب از مقیاس های

بیشتر یا هر دو رنگی خطی از آنها در آن دو خاصیت زیر باشد:

۱- چون در صورت برازش یک مدل صحیح $E(e_i) = 0$ است، لذا باید نمودار پراکنش فوق یک طرح

تصادف حول محور α ها را نشان دهد.

۲- یک تفسیر بزرگ ثابت حول محور α ها وجود داشته باشد. (چون در این مدل هیچ ثابتی نیست)
بنابراین می توان گفت که اگر در نمودارها برآینش فوق هرگونه طریقی وجود داشته باشد می توان
عبر معبر بودن مدل برایش شکی نیست.
آشکارات

~~در مدل رگرسیون چندگانه اگر هر دو شرط زیر برقرار باشد، نمودار مانده ها حاوی اطلاعات مستقیم
است که برای این سؤال که کنار مدل برایش گفته~~

نمودارهای مربوط به معبر افزوده به کمک این نمودارها بطور شهودی می توان به تاثیر هر یک
از متغیرها پیش از توجه به معبر پاسخ به دل در نظر گرفتن تاثیر سایر متغیرها پیشگویی کرد.
ابتداء مدل (*) $Y = X\beta + \epsilon$, $\text{var}(\epsilon) = \sigma^2 I$
مانده Z یا $(P+1) \times 1$ در نظر بگیرد. فرض کنید می خواهیم مقدمات درودین

متغیر پیشگویی جدید Z به مدل فوق اضافه کنیم، یعنی مدل زیر را در نظر بگیریم:

$$Y = X\beta + Z\alpha + \epsilon \quad (**) \quad Z = (z_1, z_2, \dots, z_n)' \quad \alpha \in \mathbb{R}$$

بنابراین، در ابتدا تاثیر چیزی معبر جدید Z در Y را بدون در نظر گرفتن تاثیر سایر متغیر
علاوه بر Y بررسی کنیم. برای این چنین نموداری کافی است مانده ها حاصل از رگرسیون

$Y = X\beta + \epsilon$ را بدون مانده ها حاصل از مدل $Z = X\delta + \epsilon$ رگرسیون کنیم.

معادلات رگرسیونی $e_{yx} = y - \hat{y} = (I - H_x)y$ را عنوان معبر وابسته جدید مدل معبر مستقل جدید
 $e_{zx} = z - \hat{z} = (I - H_x)z$ معبر رگرسیون کنیم بطوریکه

$$A_x = X(X'X)^{-1}X'$$

بنابراین می توان گفت که نمودارهای معبر افزوده در واقع نشان دهنده تغییرات Y است
که توسط X قابل توجه نبوده و پس از حذف اثر آن توسط Z توجه پذیر است.

در توان نوشت:

$$e_{yx} = y - \hat{y} = (I - H_x) y$$

$$e_{zx} = z - \hat{z} = (I - H_x) z$$

$(I - H_x) \epsilon$

$$y = X\beta + Z\alpha + \epsilon \xrightarrow{X(I - H_x)} (I - H_x)y = (I - H_x)X\beta + (I - H_x)Z\alpha + (I - H_x)\epsilon$$

$$\Rightarrow \text{IHK} \quad e_{yx} = (X - X(X'X)^{-1}X'X)\beta + e_{zx}\alpha + \epsilon^* \quad (\epsilon^* = (I - H_x)\epsilon)$$

$$\Rightarrow e_{yx} = 0 + e_{zx}\alpha + \epsilon^* \quad (\text{مدل مورد، مستقیم افزوده})$$

$$\Rightarrow e_{yx} = e_{zx}\alpha + \epsilon^*$$

اگر $\hat{\alpha}_{AVP}$ مقدار برابر α در مدل فوق و $\hat{\alpha}_{LS}$ برابر α در مدل $y = X\beta + Z\alpha + \epsilon$ باشد، در توان ثابت خورد ϵ برابر حالت مستقیم معبران معبران

ثابت خورد یعنی مدل $y = \beta_0 + \beta_1 X + \alpha Z + \epsilon$ را بنظر گرفته و با اضافه کردن مستقیم Z به مدل مقدار α را از روش برابر α درستی رابطه فوق را تصدیق نماید.

$$y = \beta_0 + \beta_1 X + \alpha Z + \epsilon$$

حال اگر مدل (***) صحیح باشد آنگاه مقدار مستقیم افزوده باید قابل را تولید نماید بطوریکه در تمام حالت حفر با شب $\hat{\alpha}_{LS}$ که از برابر α برآید شوند.

مثال: بررسی رابطه در مورد قیمت درخت استخوان جدید در نیویورک در این مطالعه مستقیم بصورت زیر تعیین می شوند:

- y : قیمت یک شام رجب دلار
- X_1 : هزینه غذا
- X_2 : هاج دکوراسیون
- X_3 : " " سردیس
- X_4 : مستقیم تو معینی برابر شرق و برابر غرب

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 X_4 + \beta_5 X_5 + \beta_6 X_6 + \epsilon$$

برآیند Y در شکل ۶.۹ نمودار ~~افزودن~~ هر یک از متغیرهای مستقل رسم شده. در این شکل تا اثر لغت از متغیرهای مستقل در Y غایتش دارد. شمره است. برای بررسی اثر هر یک از متغیرهای مستقل بر Y نمودارهای افزودن ~~برای~~ هر یک از متغیرهای مستقل رسم شده اند. همانطور که دیدیم، هر چه از متغیر مستقل X_1 در Y تاثیر دارد، در Y تاثیر دارد. ~~برای~~ این دو توان گفت که متغیرها دارای تاثیر اند. ~~همانطور که دیدیم~~، ~~هر چه از متغیر مستقل X_1 در Y تاثیر دارد~~، ~~در Y تاثیر دارد~~. ~~برای~~ این دو توان گفت که متغیرها دارای تاثیر اند. ~~همانطور که دیدیم~~، ~~هر چه از متغیر مستقل X_1 در Y تاثیر دارد~~، ~~در Y تاثیر دارد~~.

۲.۶ تبدیلات

در این بخش تبدیلات را بررسی خواهیم نمود که به منظور زیر بردن متغیرها (اعمال و گزینند):

۱- عمده علیه بر شکل علیه خطی بودن

۲- علیه بر شکل ثابت نبودن واریانس

۱.۲.۶ استفاده از تبدیلات برای علیه بر علیه خطی بودن

فائده فصل سوم، در این بخش نیز، در روش کلی برای یافتن تبدیل مناسب معرفی خواهند شد که عبارتند از:

۱- نمودار پاسخ معکوس (Inverse Response Plot)

۲- روش باکس-کولکس

در سه موقعیت و توانیم از تبدیلات استفاده کنیم که عبارتند از:

(a) اعمال تبدیل تنها بر روی متغیر پاسخ

(b) اعمال تبدیل بر روی متغیرها، متغیرهای مستقل

(c) اعمال تبدیل بر روی متغیرها، پاسخ و مستقل

استفاده از تبدیل به روش رگرسیون معکوس تنها در استفسار واضح :

فرض کنید مدل صحیح رگرسیون به شکل زیر باشد: $Y = g(\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p)$
که در آن $g(\cdot)$ تابع نامعلوم است. مدل فوق با استفاده از تابع معکوس $g(\cdot)$ و برآورد
به راحتی به مدل رگرسیون خطی زیر تبدیل گردد:

$$g^{-1}(Y) = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p$$

معبران مثال: $Y = \exp\{\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p\} \Rightarrow \log(Y) = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p$

مثال: در این مثال خواهیم متغیر Y (Defective) را بر روی X_1 (Temp) و X_2 (Rate) و X_3 (Rate) رگرسیون نمائیم. داده ها در جدول ۴.۱ آمده است. ابتدا با مدل زیر شروع

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \epsilon \beta_3 X_3 + \epsilon$$

شکل ۴.۱۲ نمودار فاندن آن استاندارد را در مقابل هر یک متغیرهای مستقل و برآزش شده نشان
می دهد. همه این اشکال اثر به یک طرح غیر تصادفی دارند. در شکل ۴.۱۳ نمودار برآزش \hat{Y}
در مقابل Y رسم شده که با استفاده از آن می توان به صحت یک خط راست به این نقاط
پی برد. بنابراین مدل خطی برای داده ها مناسب نبوده و باید به دنبال یک تبدیل مناسب باشیم.
بدین منظور از روش رسم پاسخ معکوس استفاده می نمائیم. فرض می کنیم مدل مناسب صورت زیر باشد:

$$Y = g(\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \epsilon) \quad \text{یا} \quad g^{-1}(Y) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3$$

با رسم $\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X_1 + \hat{\beta}_2 X_2 + \hat{\beta}_3 X_3$ در مقابل Y در شکل ۴.۱۴ می بینیم که تبدیل مناسب
 $g^{-1}(Y) = Y^{1/4.4}$. بنابراین مدل مناسب برای این داده ها عبارتست از:

$$Y^{1/4.4} = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \epsilon$$

اعمال تبدیل به روش متغیر پاسخ با استفاده از تکنیک روش بکس-کاکس و همانند گذشته، به کمک برنامه خانوار
می شود یافته تبدیلات توانی بر روی Y بصورت زیر تعریف خواهند شد:

$$\Psi_M(\gamma, \lambda) = \Psi_S(\gamma, \lambda) \cdot g_M(\gamma)^{1-\lambda} = \begin{cases} g_M(\gamma)^{1-\lambda} \cdot \frac{\gamma^\lambda - 1}{\lambda} & \lambda \neq 0 \\ g_M(\gamma) \log(\gamma) & \lambda = 0 \end{cases}$$

این روش بر این اساس است که در واقع برابر برخی مقادیر $\Psi_M(\gamma, \lambda)$ دارای توزیع نرمال تقریبی است. این روش مقدار کهنه را را با ^{مکانزیم} سازی تابع در ستمای برآورد می نماید.

مثال: بررسی مثال گذشته با زودگذر بالکس - کاکس: در شکل 4.15، تابع در ستمای $\Psi_M(\gamma, \lambda)$ در مقابل λ رسم شده است. بر این اساس، مقدار مکانزیم کهنه را برابر 1.45 خواهد بود که به جواب روش خودار معکوس بسیار نزدیک است. در ادامه، نمودار پراکنش λ در مقابل λ به متغیر مستقل مربوط رسم شده است و مث هده می شود که رابطه آنها به ^{بسیار} رابطه λ به λ نزدیک تر شده است. در شکل 4.17، نمودار پراکنش مانده λ است که در مقابل متغیر λ مستقل و مقادیر برابرش λ رسم شده که همه آنها به مناسب بودن مدل جدید تأکید دارند. در شکل 4.18، نمودار پراکنش λ در مقابل λ رسم شده است. توجه به اینده هر چه خط پرازش شده به این نقاط به نفع اول نزدیک تر شده، ^{نزدیک تر است} لذا می توان به شکی در مدل پرازش شده بود. خروجی مربوط به مدل جدید در صفحه 175 است. با استفاده از این خروجی می توان فهمید که متغیر مستقل X_3 تأثیر معنی داری بر مدل تغییرات λ نداشته و همین آنگ برابر صفر است. این مطلب با رسم نمودار λ متغیر افزوده در شکل 4.20 به وضوح فهمیده می شود (برابر اطمینان بیشتر زیرا X_3 با λ متغیر مستقل برابر است نمودار λ پراکنش اولیه رابطه خوبی دارد و لذا بهتر است اثر خروجی آن بر متغیر وابسته به λ سنجیده گردد. همچنین در صورت وجود همبستگی بین متغیر λ مستقل، آنگاه t -استدات مربوط به معنی داری ضرایب آنها (در مقابل آماره t) بیشتر می شود.

انجام تبدیلات بطور تمام بر روی متغیرهای وابسته و مستقل:

در حالتی که توزیع صد متغیرهای مستقل و پاسخ همگی دارای چولگی بوده و نیاز به تبدیل برابر نهادهای چند غیرمستقل را در نظر می‌گیریم، می‌توانیم از دو رویکرد زیر استفاده نماییم.

رویکرد ۱: این روش ترکیبی از روش بالکس-کاکس چندمتغیره و روش پاسخ معکوس است. در این روش ابتدا با استفاده از روش بالکس-کاکس چندمتغیره، تبدیل مناسب را برای متغیرهای مستقل یافته و سپس به کمک روش پاسخ معکوس تبدیل مناسب را برای متغیرهای پاسخ یافته.

$$(x_1, x_2, \dots, x_p) \xrightarrow[\text{متغیره}]{\text{بالکس-کاکس چند}} \psi_M(x_1, \lambda_{x_1}), \psi_M(x_2, \lambda_{x_2}), \dots, \psi_M(x_p, \lambda_{x_p})$$

$$y = \beta_0 + \beta_1 \psi_M(x_1, \lambda_{x_1}) + \dots + \beta_p \psi_M(x_p, \lambda_{x_p})$$

$$\text{plot}(y, \hat{y}) \rightarrow \text{یافتن بهترین مدل پاسخ} \quad (\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots + \hat{\beta}_p x_p)$$

$$y^{\text{تبدیل}} = \beta_0 + \beta_1 \psi_M(x_1, \lambda_{x_1}) + \dots + \beta_p \psi_M(x_p, \lambda_{x_p}) + \epsilon$$

رویکرد ۲: در این روش با استفاده از تبدیل بالکس-کاکس چندمتغیره، تبدیل مناسب را برای y, x_1, \dots, x_p را بدست می‌آوریم.

مثال: (سود مجله) یک تحلیلگر علاقه مندی هم ارتباط بین سود حاصل از فروش یک مجله

و آگهی‌ها آن است. او داده‌های زیر را برای بررسی سود مجله در ایالات متحده

جمع و متغیرهای وابسته و مستقل را به شرح زیر تعریف نموده است:

y : سود حاصل از آگهی

x_1 : تعداد صفحاتی که در آن آگهی (تبلیغات) درج شود

x_2 : رتبه حاصل از آگهی مشتریان

x_3 : رتبه حاصل از رتبه‌ها

در شکل ۶.۲۱ نمودار پراکنش متغیرها رسم شده که نتایج کلی آن نشان دهنده وجود چگالی در کلیه متغیر است. همچنین به نظر می آید که بین متغیرهای مستقل رابطه خطی وجود داشته باشد بنابراین می توان به اعمال تبدیل در متغیرهای مستقل وابسته روی آن نمود.

اندازه استاندارد از سه اعمال تبدیلیات روی متغیرها استوار از روش در این است. برای این اساس ابتدا تبدیل مناسب برای متغیرهای مستقل را با استفاده از بکس-کالکس چند متغیر بدست آوریم. بنا بر فرضی صفحه ۱۷۷ می توان گفت که λ مربوط به هر سه متغیر مستقل تقریباً برابر همگراست. پس، تبدیل مناسب تبدیل لگاریتمی است. بنابراین:

$$Y = \beta_0 + \beta_1 \log(X_1) + \beta_2 \log(X_2) + \beta_3 \log(X_3) + \epsilon \quad (1)$$

در شکل ۶.۲۲ رسم نمودار معکوس λ در مقابل \hat{Y} درجیم که تبدیل مناسب برای λ به ازای $\log Y = -0.20$ حاصل می شود. بنابراین مدل نهایی بصورت زیر خواهد بود:

$$Y^{-0.20} = \beta_0 + \beta_1 \log(X_1) + \beta_2 \log(X_2) + \beta_3 \log(X_3) + \epsilon \quad (2)$$

اما با توجه به شکل $\lambda = 0$ نیز برابر مقادیر کوچک و متوسط متغیر λ ، برازش خوبی به داده ها دارد و چون همچنین متغیر مستقل وابسته یک است، لذا تبدیل جانشین دیگر و روانه $\log(Y)$ باشد و مدل بصورت زیر در نظر گرفته شود:

$$\log(Y) = \beta_0 + \beta_1 \log(X_1) + \dots + \beta_3 \log(X_3) + \epsilon \quad (3)$$

اما به هر حال برابر هر کدام از دو تبدیل فوق گفته است آنرا که جدول آن نیز در این فایل قابل قبولتری را در هر دو عنوان مدل های در نظر بگیریم.

با اعمال تبدیل بر روی متغیرها با استفاده از روش درجیم برای λ صفحه ۱۷۹ ، می توان نتیجه گرفت که روش بکس-کالکس برابر کلیه متغیرهای تبدیل لگاریتمی می تواند تبدیل مناسب معرفی می نماید. در شکل ۶.۲۳ ، نمودار پراکنش مارتسی برای هر سه متغیر رسم شده. همانطور که از این شکل برآید، متغیرهای تبدیل یافته رابطه خطی معنی دارتری با هم دارند.

در شکل ۶۰۲۳، نمودار ~~عیب شناسی~~ ~~برای~~ ~~مدل~~ ~~لگاریتمی~~ رسم شده است که

گهی که نه صفتا طریقی بقا در آن هستند. بنابراین مدل فوق مدلی مناسب برای داده ها نیست.

همچنین در شکل ۶۰۲۵، نمودار ^{جدید} ~~برای~~ ~~مدل~~ ~~لگاریتمی~~ تبدیل یافته رسم شده و چون

نقاط تقریباً بر روی نیمه از ربع اول و سوم واقعند، لذا مدل برازش کافی به داده ها دارد.

در شکل ۶۰۲۶، نمودار ~~عیب شناسی~~ ~~برای~~ ~~مدل~~ ~~جدید~~ رسم شده است که از مدل (۳) که گهی

ملاقات بر آن مناسب بودن مدل مورد بررسی دارند. البته نمودار پایین دست چپ این شکل که نشان دهنده

که وار ^{خطا} ~~ن~~ ~~افزایش~~ ~~در~~ ~~سین~~ ~~کاهش~~ دارد. همچنین در نمودار پایین دست راست نقطه چین عمودی

شکل نه صفتا نیز لوریج ها یعنی $\frac{2(P+1)}{n} = 1.09$ و خط ^{افقی} ~~عمودی~~ ~~ن~~ ~~نه~~ ~~صفتا~~ مرکزیت مانده کار استناد دارد

است.

در شکل ۶۰۲۷، نمودار ~~عیب شناسی~~ ~~برای~~ ~~مدل~~ ~~جدید~~ رسم شده که بر اساس آن سوال ^{فهرده} ~~که~~ ~~مستخرج~~

$\log(x_c)$ تأثیر معنی داری روی دور مدل (۳) ندارد و بهترین مدل حذف گردد. این مطلب

با استناد به بار خروجی صفتا ۱۸۴ نیز کاملاً آشکار است.

برای در رسید تعبیرات مستخرج است. برابر ~~تفسیر~~ ~~تقریباً~~ ~~از~~ ~~مستخرج~~ ~~مستعمل~~ : مدل رگرسیون زیر

رابطه تعریف شده :
$$\log(y) = \beta_0 + \beta_1 \log(x_1) + \beta_2 \log(x_2) + \epsilon$$

باید بودن مستخرج x_2 سوال نیست :

$$\beta_2 = \frac{\Delta \log(y)}{\Delta x_2} = \frac{\log(y_2) - \log(y_1)}{\Delta x_2} = \frac{\log(y_2/y_1)}{\Delta x_2} = \frac{y_2/y_1 - 1}{\Delta x_2} = \frac{\% \Delta y}{100 \Delta x_2}$$

$$\Rightarrow \% \Delta y = 100 \Delta x_2 \times \beta_2 \quad \Delta x_2 = 1 \quad \Rightarrow \boxed{\% \Delta y = 100 \beta_2}$$

مثال: برای دور تیراژ یک روزنامه در یکشنبه ۱۸۴

در این مثال بر اساس ۱۸۹ مشاهده در ایالات متحده آمریکا مستخرج زیر تعریف شده است:

y : لگاریتم تیراژ در یکشنبه

x_1 : لگاریتم تیراژ در سایر روزهای هفته

x_2 : مستخرج در حالی که سفر در یک خط هوایی به شهر نیویارک در روز یکشنبه

و شورو یک بار حالتی است که در آن شهر روزنامه دیگری که حالتی خلاصه دارد نیز علاوه بر روزنامه اصلی چاپ و شورو.

در شکل ۴.۲۸، نمودار پیکانش γ در مقابل X_1 رسم شده است (برابر مدل $\hat{y} = 10.0 + 1.1X_1$ و $R^2 = 0.97$)
 در شکل ۴.۲۹ نمودار فاندان استناد در مقابل متغیر مستقل و برابرش شده رسم شده است که همین نشانگر یک طرح تصادفی هستند. شکل ۴.۳۰ نمودار پیکانش γ را در مقابل $\hat{\gamma}$ نشان دهد که با استفاده از آن سوال به شیوهی برابرش می برد. در شکل ۴.۳۱ نمودار بار مربوط به عیب ای دیگر بدون رسم شده اند که همه آنها اعتبار مدل نگارشی را تایید و نمایند. خروجی R مربوط به مدل فوق در صفحه ۱۸۹ آمده که برابر با آن سوال به معنی داری هر دو متغیر مستقل X_1 و X_2 می برد. برابر با آن سوال گفت:

- (الف) ثابت بودن متغیر X_2 ، افزایش متغیر X_1 به ازای X_2 برابر است با $\frac{\partial \hat{y}}{\partial X_1}$ یا $\frac{1.1}{1.0}$ یا 1.1
 (ب) X_1 ، متغیر X_2 به ازای X_1 در هر مقدار X_2 مقدار $\frac{\partial \hat{y}}{\partial X_2}$ یا 0.1 کاهش می یابد.

در شکل ۴.۳۲ ، نمودار متغیر افزودن برابر مدل فوق رسم شده اند که با استفاده از آن سوال به معنی داری اثر هر یک از متغیرها مستقل بر این مدل می برد.

در این مورد ما قادر خواهیم بود که مقدار متغیر γ را به ازای تعداد لغات تراز هستند $1-1000$ بصورت زیر برآورد و فاصله اطمینان ۹۵٪ بدست آوریم:

$$X_2 = 1 \rightarrow \hat{\gamma} = 12,028 \quad (11,721, 12,355)$$

$$X_2 = 0 \rightarrow \hat{\gamma} = 12,559 \quad (12,280, 12,838)$$

$$X_2 = 1 \rightarrow \text{تعداد تراز تپسته} = \exp(12,028)$$

$$X_2 = 0 \rightarrow \text{ " " " } = \exp(12,559)$$

همچون چندگانه اگر بین دریا چند متغیر مستقل در مدل رگرسیون چندگانه رابطه خطی وجود داشته باشد، گوئیم مدل رگرسیون را از همخطی است. همخطی کامل زمانی رخ دهد که یکی از متغیرهای مستقل تابع دقیق از بقیه متغیر مستقل دیگر باشد و همخطی ناقص زمانی خواهد رخ داد که این تابع تقریبی باشد. در عمل همخطی کامل رخ نخواهد داد زیرا اگر یکی از متغیرهای مستقل تابعی دقیق از سایر متغیرهای دیگر باشد، خطای اندازه گیری بسیار و شونده این رابطه دقیق به یک رابطه تقریبی تبدیل گردد. حال چنانچه همخطی تقریبی اتفاق افتد، آنگاه $|X'X| \rightarrow 0$ و لذا ماتریس $(X'X)^{-1}$ وجود نخواهد داشت.

منابع اصلی همخطی:

- ۱- انتخاب روش گردآوری داده‌ها
- ۲- گذاشتن متغیری روی ^{مدل} X جامع
- ۳- شناسایی X مدل

۴- انتخاب مدل که تعداد متغیرهای مستقل آن بیشتر از تعداد مشاهده باشد.

در حالت کلی با زیاد شدن متغیرهای بیشتر، اهمیت بوجود آمدن همخطی بین متغیرهای مستقل افزایش می‌یابد.

اثرات همخطی: همانطور که گفته شد، اگر همخطی وجود داشته باشد، $|X'X| \rightarrow 0$ شده و لذا دارایی‌های ماتریس $(X'X)^{-1}$ خیلی بزرگ خواهند شد و چون $\text{Var}(\hat{\beta}) = \sigma^2 (X'X)^{-1}$ ، لذا می‌توان نتیجه گرفت که وقت برآورد کمترین مربعات کاهش یافته و همچنین عوامل اطمینان بدست آمده برابر پارامترها بین خواهند شد. گاه اوقات مسخ است حتی علامت پارامترهای β_1, \dots, β_p به اشتباه برآورد گردد. همچنین مقایسه آماره‌های t مربوط به ضرایب رگرسیون بسیار کوچک شده و غیرمعنی‌دار به نظر می‌رسد زیرا $t_{\beta_j} = \frac{\hat{\beta}_j}{\sqrt{\sigma^2 c_{jj}}}$ که در آن t_1, t_2, \dots, t_p (توان آمین) را که در صورت مقصود اصلی ماتریس $(X'X)^{-1}$ است.

چگونه می توان به وجود همخطی پی برد؟
 کانه اوقات آماره F جدول آماره های متنوع
 معنی دار است می تواند از آماره t مربوط به ضرایب رگرسیون معنی دار نشود
 در چنین مواردی می توان به وجود همخطی پی برد و می بایستی از مجموعه داره های که همخطی معنی
 داری دارند این رفتار را شرح می دهند و لذا این معیار چندان سودمند نمی باشد. در چنین
 مواردی از معیار زیر می توان استفاده نمود. هر آنگاه که:

Condition number
 ۱- عدد شرطی: این آماره بر اساس مقادیر ویژه ماتریس $(X'X)$ بصورت زیر تعریف
 می شود:

$$K(X'X) = \sqrt{\frac{\max \lambda_i}{\min \lambda_i}} \quad i=1, \dots, p$$

که در آن $\lambda_1, \dots, \lambda_p$ مقادیر ویژه ماتریس $(X'X)$ هستند.

۲- عامل تورم واریانس (Variance Inflation Factor) عامل تورم واریانس مربوط به متغیر مستقل
 نام j را VIF_j نشان دارد و بصورت زیر تعریف می کنیم:

$$VIF_j = \frac{1}{1 - R_j^2}$$

که در آن R_j^2 ضریب تعیین در مدل رگرسیون X_j در برابر سایر متغیرهای مستقل است.

هرگاه $VIF_j > 5$ (یا $R_j^2 > 0.8$) آنگاه می توان به وجود همخطی پی برد.

همانطور که در رگرسیون $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \epsilon$ در مدل خطی در متغیرها

$$S_{X_j X_j} = \sum_{i=1}^n (X_{ij} - \bar{X}_j)^2$$

$$\text{Var}(\hat{\beta}_j) = \frac{1}{1 - r_{j1}^2} \cdot \frac{\sigma^2}{S_{X_j X_j}}$$

بنابراین می توان نوشت:

$$\text{Var}(\hat{\beta}_j) \rightarrow \infty \Rightarrow VIF_j \rightarrow \infty \Rightarrow \text{Var}(\hat{\beta}_j) \rightarrow \infty$$

وقت بزرگوارتر کمترین مربعات کم می شود.

۲۴
 رابطه اخیر در حالت کلی نیز درست بود و در آن نشان دادیم: $y = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p + \epsilon$

$$\text{Var}(\hat{\beta}_j) = \frac{1}{1-R_j^2} \cdot \frac{\sigma^2}{S_{X_j X_j}} \quad j=1, \dots, p$$

$$\Rightarrow \text{var}(\hat{\beta}_j) = \text{VIF}_j \cdot \frac{\sigma^2}{S_{X_j X_j}}$$

ص ۹۵

مثال: (سافت پیل)

در سافت پیل، متغیرهای زیر قابل بررسی هستند:

۱: زمان سافت بر حسب روز (Time)

X_1 : مساحت پیل بر حسب فوت مربع (DArea)

X_2 : هزینه سافت بر حسب دلار (Cost)

X_3 : تعداد چراغان پیل (DWgs)

X_4 : طول پیل (Length)

X_5 : تعداد طاق‌ها (Spans)

ابتدا به کمک روش بکس-کاکس چند متغیره، تبدیل مناسب جهت تبدیل داده‌ها به یک توزیع نرمال چند متغیره، و ای در رابطه اصلی بین آنها را بدست می‌آوریم. سوال اصلی این است که آیا متغیرها برابر هستند در تمام گرفت و درستی مدل زیر حاصل شود.

$$\text{Log}(y) = \beta_0 + \beta_1 \text{Log}(X_1) + \dots + \beta_5 \text{Log}(X_5) + \epsilon$$

شکل ۶.۳۹ و ۶.۴۰: بترتیب فائز نمودار پیرائنت متغیر را قبل و بعد از تبدیل نشان داده‌اند.
 شکل ۶.۴۱: نمودار فاندمان است که در مقابل متغیرها را بسته و مستقل برابر مدل تبدیل یافته را نشان می‌دهد. (همه نمودارهای نشان داده شده برای تعدادن هستند).

شکل ۶.۴۲: نمودار $\text{Log}(y)$ در مقابل $\text{Log}(X_1)$ و در شکل ۶.۴۳ نمودارهای مربوط به یکسای رگرسیون رسم شده‌اند که کجی دلالت بر متغیر بودن مدل فوق دارند. در

اما هنگامی که به جدول آنالیز واریانس برای خودی می بینیم که با وجود آنکه آماره F در سطح
 بالایی معنی دار است، آماره t مربوط به متغیر X_3 معنی دار بوده و سایر متغیرها
 در این اثر معنی داری ندارند. علاوه بر این علامت مربوط به

ضرایب برابر شده، متغیرهای مستقل X_1 (ساعت) و X_2 (طول پل) معنای معنی
 هر چه سطح (طول) یک پل افزایش یابد زمان ساخت آن کمتر شود که این جمله
 منطقی است. همچنین براساس شکل (۶.۴۵) (عنوان آن متغیر افزوده) آماره

متغیر مستقل X_3 در این اثر معنی دار بر مبنای جدول بوده و سایر متغیرها فاقد اثر معنی دار بر مهند
 هرگاه رویا چند متغیر مستقل با همگی بالا وارد یک مدل رگرسیون شوند، آنها بطور
 مؤثری اطلاعات مشابه را در مورد متغیر وابسته حمل کرده و این باعث می شود که روش

کمترین مربعات نتواند اثر جزئی هر یک از آنها را در متغیر پاسخ تشخیص دهد. در چنین
 مواقعی آماره F جدول آنالیز واریانس بزرگ شده و سطح معنی داری آن بسیار
 افزایش می یابد در حالی که تعداد کم از ~~آماره کل مربوط به~~ ضرایب معنی دار ~~شوند~~

مشکل دیگر ایجاد شده در چنین مواقعی، بر آورد علامت ضرایب رگرسیون می باشد.
 ما این همبستگی متغیرهای مستقل در صفحه ۲.۲ آمده است. همانطور که دیدیم، در مورد این
 متغیرها همبستگی ~~مستقل~~ حتی بالایی وجود دارد. همچنین مقادیر VIF برای

X_1	X_2	X_3	X_4	X_5
۷,۱۱۴	۸,۱۴۸	۲,۱۴۱	۸,۱۰۱	۲,۱۸۸

"کوی انتخاب متغیر در مدل"

در این بخش روشهای برای انتخاب بهترین مدل ممکن از میان تمام مدلها قابل انتخاب، معرفی خواهد شد.

مدل دگرسیون چندگانه زیر را در نظر بگیرید:
$$Y = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p + \epsilon$$

روشهای انتخاب متغیرها به انتخاب بهترین زیر مجموعه ممکن از متغیرهای مستقل است. در ادامه معیارهایی جهت ارزیابی زیر مجموعه‌های از متغیرهای مستقل ارائه خواهد شد:

۱- معیار R^2_{adj} : همانطور که در نامش اضافه نمودن متغیرهای مستقل بی ربط اغلب باعث افزایش R^2 و شونز وی این اتفاق در مورد R^2 تصحیح شد، که بصورت زیر تعریف و توجیه خواهد شد:

$$R^2_{adj} = 1 - \frac{SSE / (n - p - 1)}{S_{yy} / (n - 1)}$$

در آن ثابت نمود که در صورتی افزایش R^2_{adj} با افزودن یک متغیر به مدل اتفاق خواهد افتاد هرگاه F آماره F چیزی مربوط به آن متغیر مستقل ازین بیشتر باشد. در عمل تعداد از متغیرهای مستقل طوری انتخاب و شود که منجر به بیشترین آماره R^2_{adj} شوند. همان تان دار کالزیم R^2_{adj} برای متغیرهای مستقلی اتفاق خواهد افتاد که دارای کمترین SSE باشند.

البته چنانچه افزایش یک متغیر منجر به بهبود مدل باعث افزایش چیزی در R^2_{adj} شود برای تفسیر بهتر نتایج رسانی بیشتر، کفهرات آن متغیرها در مدل‌های در نظر بگیریم. بعنوان مثال فرض کنید:

$P = 9 \rightarrow R^2_{adj} = 0.1791$ $P = 10 \rightarrow R^2_{adj} = 0.1792$

$P = 8 \rightarrow R^2_{adj} = 0.1541$

در مدلها فوق هر چند $P = 1$ متغیر مستقل در این بیشترین مقدار R^2_{adj} است ولی کفهرات که $P = 9$ متغیر در نظر گرفته شود زیرا تفسیر چندانی در R^2_{adj} ایجاد نمی‌کند.

سه معیار زیر تنها در صورتی قابل استناد است که توزیع معیارها یکسان باشد و البته نزنال باشد در صورت نزنال بودن توزیع آنها تابع درستی برابر است:

$$\log L(\beta_0, \dots, \beta_p, \sigma^2 | y) = -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{i1} - \dots - \beta_p x_{ip})^2$$

$$\rightarrow \log L(\hat{\beta}_0, \dots, \hat{\beta}_p, \hat{\sigma}^2 | y) = -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log(\hat{\sigma}^2) - \frac{1}{2\hat{\sigma}^2} SSE$$

$$\rightarrow \hat{\sigma}^2_{MLE} = \frac{SSE}{n} \quad , \quad \hat{\sigma}^2_{LS} = S^2 = \frac{SSE}{n-p-1}$$

$$\rightarrow \log L(\hat{\beta}_0, \dots, \hat{\beta}_p, \hat{\sigma}^2 | y) = -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log\left(\frac{SSE}{n}\right) - \frac{n}{2}$$

۲- معیار AIC (Akaike's Information Criterion)

این معیار بر اساس گزینش تابع درستی تعریف شده و سئوی برازش را اندازه میگیرد. AIC همیاری جهت اندازه گیری اطلاعات بود و بصورت زیر تعریف میشود:

$$AIC = -2 \left[\log L(\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p, \hat{\sigma}^2 | y) - (p+1) \right]$$

$$= -2 \left[-\frac{n}{2} \log(2\pi) - \frac{n}{2} \log\left(\frac{SSE}{n}\right) - \frac{n}{2} - (p+1) \right]$$

$$= n \log\left(\frac{SSE}{n}\right) + 2p + \text{other terms}$$

بطوریکه سایر جملات به نحو برازش مدل بستگی نداشته و برابر همه مدلهاست. بنابراین AIC در نرم افزار R بصورت زیر محاسبه میشود:

$$AIC = n \log\left(\frac{SSE}{n}\right) + 2p$$

لازم به توضیح است که بهترین مدل مدلی است که دارای کمترین AIC باشد.

۳- معیار AIC تصحیح شده (AIC_c): اگر حجم نمونه کوچک و یا تعداد پارامترها

مدل نسبتاً زیاد باشد، کجراهی از AIC_c استوار شود. اگر $\frac{n}{p+1} \leq 4$ باشد، AIC_c باید بکاربرد شود. این معیار بصورت زیر محاسبه میشود:

$$AIC_c = AIC + \frac{2(p+1)(p+2)}{n-p-1}$$

$$AIC_c \xrightarrow{n \rightarrow \infty} AIC$$

۴- معیار BIC (Bayesian Information Criterion)

معیار اطلاع بنر بصورت زیر محاسبه می شود:

$$BIC = -2 \log L(\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p, \hat{\sigma}^2 | y) + (p+1) \log(n)$$

مشابه معیار AIC، هرچه BIC کوچکتر باشد مدل برازش بهتری به داده دارد.

۵- معیار Cp مالوس: این معیار تخمین برآورد متوسط مالوس پیشفاده و بصورت زیر تعریف شده:

$$C_p^* = \frac{SSE_{p^*}}{s^2} - \frac{(n-2p^*)}{n+2(p+1)}$$

که در آن SSE_{p^*} خطای p^* پارامتر از جمله β_0 و s^2 برآورد σ^2 با MSE بیشترین تعداد متغیرها را مستلزم است. حال اگر مدلی با p^* پارامتر کیفیت کند در اینصورت $E(SSE_{p^*}) = (n-p^*)\sigma^2$ و چون باین فرض $E(s^2) = \sigma^2$ است لذا می توان نوشت:

$$E(C_p^*) \cong \frac{(n-p^*)\sigma^2}{\sigma^2} - (n-2p^*) = p^*$$

بنابراین می توان نتیجه گرفت هرگاه مقدار C_p^* به p^* نزدیک شود، مدل مناسب خواهد بود. تصمیم گیری در مورد انتخاب کالسیون همواره زیر مجموعه آن مستلزم است. این معیار را می توان به این نام تصمیم گیری در مورد متغیرها مفید در مدل در راه کلی زیر که کاملاً متفاوت وجود دارد:

۱- همه دیگر معیارها مستلزم است؛ این روش بر اساس برازش همه 2^p مدل درگیر کردن مستلزم است به متغیرها و شناسایی زیر مجموعه آن از متغیرها مستلزم است که معیار برازش را ماکزیمم و یا معیار اطلاع را مینیمم نماید. البته باید توجه داشت که معیارها متفاوت مستلزم است متغیرها نتایج متفاوتی گردند و هیچ معیاری به تنهایی همیشه بهترین نخواهد بود. لذا باید بیشتر از یک معیار را در نظر بگیریم.

مثال: عددی که سافت پل بر آن گزینیم. ابتدا به کلی شغل تمام متغیرها را در نظر می‌گیریم:

$$\text{Log}(Y) = \beta_0 + \beta_1 \text{Log}(X_1) + \dots + \beta_5 \text{Log}(X_5) + \epsilon$$

خود R مربوط به مدل فوق در صفحه ۲۳۴ کتاب آمه‌ات. همانطور که مشاهده شود با وجود معنی دار بودن آماره F جدول آنالیز واریانس، غالب آماره دل + جزئی تغییرات ضریب مربوط به متغیر X_3 معنی دار نمی‌باشند. بنابراین تحقیق در مورد انتخاب متغیرهای بیشتر را، ما گزینیم سازی R^2_{adj} شروع می‌کنیم. در شکل ۷.۱ نمودار R^2_{adj} در مقابل تعداد متغیرهای مستقل رسم شده است. عنوان مثال، R^2_{adj} متغیر X_3 متغیر X_2 و X_5 را شامل می‌کند. در جدول ۷.۱ مقادیر R^2_{adj} ، AIC_c ، AIC و BIC برای زیرمجموعه دل بسته با حجم حال یک، ۵ آورده شده است. طبق این جدول، در معیار R^2_{adj} ، AIC_c و AIC بهترین مجموعه ممکن را شامل X_3 ، X_5 و X_2 پیشنهاد می‌دهد، ولی در معیار BIC و AIC_c بهترین مجموعه ممکن را شامل X_3 و X_5 در نظر می‌گیرند. با مقایسه آماره دل این جدول و تفاوت بسیار جزئی آماره دل R^2_{adj} و AIC برای دو سه متغیر مستقل بهترین انتخاب ممکن مدلی شامل X_3 و X_5 به صورت زیر است که خودی مربوط به آن در صفحه ۳۳۶ آمده است:

$$\text{Log}(Y) = \beta_0 + \beta_3 \text{Log}(X_3) + \beta_5 \text{Log}(X_5) + \epsilon$$

۲- گزینش به روش گام به گام (Stepwise Selection)

این روش در واقع جنبه اصلاح شده روش رگرسیونی پیشرو است که به ما اجازه می‌دهد در هر مرحله متغیرهای لحاظ شده در مدل در مراحل قبل را دوباره امتحان کنیم. متغیری که در مرحله قبل در مدل وارد شده، ممکن است در مرحله بعدی بخاطر همبستگی آن با سایر متغیرهای مستقل در اکثر زائده به نظر برسد. برای بررسی این موضوع در هر مرحله یک آزمون F جزئی برای هر متغیر که در مدل است انجام می‌شود و لو آنکه این متغیر جدیدترین متغیر

۱۷
 دارد. به مدل رگرسیون با کمینه و صرف نظر از اینکه چه وقت به مدل وارد شده است. در هر مدل متغیر با کوچکترین F چیزی بی معنی (در صورت وجود) حذف شود و مدل با متغیرهای باقیمانده دوباره برازش می شود و F های چیزی محاسبه می گردند و به طور مشابه امتحان می گردند و این کار ادامه پیدا می کند. این فرآیند تا زمانی ادامه می یابد که متغیرهای بیشتری برآورد شوند و یا از مدل حذف می گردند. بنابراین روش گام به گام

حد اکثر
$$1 + 2 + \dots + (p-1) + p = \frac{p(p+1)}{2}$$
 مدل از میان تمام 2^p مدل

ممکن را بررسی می نماید. لذا، نمی توان گفت که روش گام به گام لزوماً مدلی را ارائه می نماید که بهترین مدل از لحاظ کمترین معیار در اطلاع گفته شده در قسمت قبل است. به عبارت دیگر نتایج روش گام به گام با نتایج معیار در اطلاعاتی در R^2 ممکن است یکسان نباشد. با پیارسازی این روش در داده های مثال قبل، روش پیشرو و پسرو هر دو همگرا به بهترین مدل ممکن را می یابند. متغیرهای X_3 ، X_5 و X_2 معرفی خواهند نمود.