

دانشور

تعیین حجم نمونه در مطالعات رگرسیون خطی معمولی

نویسندها: دکتر اتوشیروان کاظم نژاد* و سید مهدی سادات هاشمی*

* استادیار آمار زیستی دانشگاه تربیت مدرس

** دانشجوی دکترا آمار زیستی دانشگاه تربیت مدرس

چکیده

در تحقیقات گوناگون، نمونه‌گیری و تعیین حجم مناسب برای دستیابی به نتایج توانمند آماری از اهمیت ویژه‌ای برخوردار است. در حیطه کار با مدل‌های رگرسیون خطی معمولی مسئله مهم این است که برای بررسی رابطه معنادار بین یک متغیر وابسته با چندین متغیر مستقل باید به چه صورت عمل کرد؛ زیرا نمی‌توان با تعیین حجم نمونه بر مبنای رابطه متغیر وابسته با تنها یک متغیر مستقل، امید داشت که این حجم تعیین شده توان آزمون‌ها برای سایر متغیرها را نیز تضمین کند. بنابراین باید با تعدیل‌های مناسب بر روی حجم نمونه اولیه، شرایط مناسب را ایجاد کرد. در این مقاله سعی شده بدون پرداختن به مسائل استنبط آماری و مشکلاتی که در صورت کافی نبودن حجم نمونه پیش خواهد آمد، روش ساده‌ای برای تعیین حجم نمونه در مطالعات رگرسیونی چندگانه ارائه گردد.

واژه‌های کلیدی: رگرسیون چندگانه، ضریب همبستگی، حجم نمونه

دوماهنامه علمی-پژوهشی
دانشگاه شاهد
سال هشتم-شماره ۳۲
اردیبهشت ۱۳۸۰

پیروی می‌کرده و تعداد کل ساعتی که در طول این مدت به ورزش می‌پرداخته را نشان دهد و از طرف دیگر، تمایل دارد تا میزان و جهت این ارتباط را نیز تعیین و یک مدل ریاضی ارائه کند تا بهوسیله آن، توانایی پیش‌بینی وزن بیماران را در آینده برحسب این دو متغیر داشته باشد. در این‌گونه موارد به طور طبیعی با مدل‌های رگرسیونی که ساده‌ترین آن‌ها مدل رگرسیون خطی معمولی است، مواجهیم. به طورکلی در این مدل‌ها، هدف، پیش‌بینی یک

مقدمه

در مطالعات گوناگون مواردی پیش می‌آیند که محقق به ارتباط بین دو یا چند متغیر علاقمند است و همچنین تمایل دارد که در صورت وجود روابط معنادار بین متغیرها، مدلی ریاضی برای پیش‌بینی یک یا چند متغیر برحسب سایر آن‌ها به دست آورد. برای مثال فرض کنید یک پزشک علاقه‌مند است ارتباط معنادار بین وزن یک بیمار با تعداد هفته‌هایی که از یک رژیم غذایی خاص

پراکنش نگار تجارت قبلی این را نشان داده باشد) هدف عبارت خواهد بود از برآورد کردن ضریب β_1 و آزمون کردن آن در معادله زیر:

$$Y = \beta_0 + \beta_1 X + e \quad (2)$$

در حالت کلی برای نشان دادن رابطه خطی بین دو متغیر که فرض می‌شود دارای توزیع نرمال توانم هستند، از ضریب همبستگی پیرسون (ρ) استفاده می‌شود که عبارت است از:

$$\rho_{xy} = \frac{COV(X, Y)}{\sqrt{Var(X)}\sqrt{Var(Y)}} \quad (3)$$

و می‌توان نشان داد که رابطه‌ای بین β و ρ_{xy} برقرار است:

$$\rho_{xy} = \frac{\beta_1 \sigma_x}{\sigma_y} \quad (4)$$

که در آن $(\sigma_x^2 = Var(X))$ و $(\sigma_y^2 = Var(y))$ و بنابراین $\rho_{xy} = \frac{\beta_1 \sigma_x}{\sigma_y}$ اگر و فقط اگر $\beta_1 = 0$. چنانچه X و Y استاندارد شده باشند (یعنی $\mu_x = \mu_y = 0$ و $\sigma_x = \sigma_y = 1$) آنگاه آزمون فرض موردنیاز برای این آزمون برآبرو یکدیگر خواهد بود [1]. حال اگر $\rho_{xy} = 0$ برآورده از ضریب همبستگی بین X و Y باشد، فرمول حجم نمونه برای آزمون کردن فرض $H_1: \rho_{xy} \neq 0$ مطابق رابطه زیر است [2]:

$$n = \frac{\left(Z_{1-\alpha/2} + Z_{1-\beta} \right)^2}{[C(r)]^2} + 3 \quad (5)$$

که در آن $C(r) = \frac{1+r}{1-r}$ که همان تبدیل فیشر، و سطح معناداری آزمون و $\alpha - \beta$ توان آن است.

ب) حالت رگرسیونی چندگانه یا بررسی ارتباط بین یک متغیر وابسته یا چند متغیر مستقل در این مرحله فرض می‌کنیم که بیش از یک متغیر تبیینی (مستقل) داشته باشیم؛ یعنی هدف بررسی ارتباط بین متغیر Y و متغیرهای X_1, X_2, \dots, X_p و پیش‌بینی Y با استفاده از مدل رگرسیونی آن بر حسب X ها باشد. به

یا چند متغیر Y (وابسته) - که فرض می‌شود دارای توزیع نرمال و پیوسته‌اند - براساس یک یا چند متغیر تبیینی (مستقل) X_1, \dots, X_p است. در این حال، معادله‌ای که Y و X را به صورت خطی به یکدیگر پیوند می‌دهد عبارت است از:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p + e \quad (1)$$

که در آن e متغیری است که خطاهای را نشان می‌دهد. معمولاً فرض می‌شود خطاهای دارای توزیع نرمال با میانگین صفر و واریانس σ_e^2 و مقادیرشان مستقل از یکدیگر هستند. اهدافی که محقق دنبال می‌کند هم ارز با معنادار بودن ضریب هر β_j (یعنی β_j) در مدل رگرسیون خطی ۱ است. به عبارت دیگر، معنادار بودن β_j نشان‌دهنده رابطه β_j با Y ، $|\beta_j|$ بیانگر میزان تأثیر X_j بر Y و علامت β_j نشان‌دهنده جهت رابطه بین β_j و Y است.

در مطالعات آماری با توجه به اهداف مورد بررسی باید نمونه‌ای را با حجم لازم از جامعه استخراج کرد تا براساس آن بتوان با خطأ و توان مناسب، استنباط‌هایی را درباره پارامترها انجام داد. از آنجاکه بررسی‌های رگرسیون نیز در حیطه مطالعات آماری قرار دارد و به طور گسترده‌تر نیز در همه زمینه‌های علمی به کار گرفته می‌شود باید این موضوع در آن‌ها مورد توجه قرار بگیرد. در این مقاله سعی می‌کنیم در ابتدا فرمولی ساده که حجم نمونه را برای نشان دادن رابطه بین یک متغیر که می‌تواند وابسته باشد و یک متغیر دیگر که می‌تواند مستقل باشد ارائه دهیم. در مرحله بعد این فرمول را برای به دست آوردن حالاتی که در آن، هدف مطالعه نشان دادن روابط معنادار بین یک متغیر وابسته و چند متغیر مستقل است، تعمیم می‌دهیم.

۲- مواد و روش‌ها

الف) حالت رگرسیونی ساده یا بررسی ارتباط بین فقط دو متغیر

کار را با دو متغیر X و Y آغاز و فرض می‌کنیم که هدف، بررسی رابطه بین این دو و همچنین در صورت وجود ارتباط، پیش‌بینی Y بر حسب X باشد. با این فرض که تغییرات Y بر حسب X خطی است (مثلاً نمودارهای



اولیه به دست آمده از رابطه ۵ برای بررسی رگرسیون
چندگانه γ بر روی X ها به صورت زیر تبدیل می‌شود [۱]:

$$\text{Eq} \quad n_p = n \times VIF^6$$

که در آن n و n_p به ترتیب عبارتند از حجم نمونه مورد نیاز برای مدل رگرسیونی با ۱ و p متغیر تبیینی، مشکلی که باقی می‌ماند این است که چگونه VIF را محاسبه کنیم. ما نحوه این محاسبه را در پیوست ۱ آورده‌ایم.

۳- کاربرد

حال با مثالی کاربردی نحوه استفاده از این فرمول را نشان می‌دهیم.

فرض کنیم پزشکی می‌خواهد رابطه بین فشار خون بیماران و میزان کلسترول خون آن‌ها را بررسی کند. با سطح معناداری ۵ درصد و توان ۹۵ درصد، حجم نمونه لازم در هر یک از موارد زیر چه تعداد بیمار خواهد بود؟
(الف) در صورتی که از مطالعات مشابه حدس زده شود که ضریب همبستگی پیرسون بین فشار خون و میزان کلسترول آن حدود 0.5 است.

(ب) از مطالعات قبلی مشخص شده باشد که بین میزان کلسترول خون با دوز مصرفی یک نوع داروی ضد فشار خون و تری‌گلیسیرید آن رابطه رگرسیونی با ضریب تعیینی $R^2 = 0.6$ وجود دارد و هدف از مطالعه نیز به دست آوردن الگوی پیش‌بینی فشار خون بر حسب دوز مصرفی دارو و میزان کلسترول خون و تری‌گلیسیرید آن باشد.

با در نظر گرفتن:

$$z_{1-\alpha/2} = 1/96, z_{1-\beta} = 1/64$$

در حالت اول با استفاده از رابطه ۵، تعداد بیماران را حداقل ۴۶ نفر و در حالت دوم با استفاده از رابطه ۶ تعداد بیماران را حداقل ۱۱۵ نفر به دست می‌آوریم. از آن جا که اساس استنباط‌هایی که در رگرسیون صورت می‌گیرد، آزمون F و آزمون t است، در صورتی که هدف، مدل‌بندی رگرسیونی فشار خون بر حسب سایر متغیرها بوده، و حجم نمونه کفايت لازم را نداشته باشد، توان آزمون‌های t و F تحت تأثیر قرار خواهد گرفت.

عبارت دیگر، فرض مورد علاقه در هر حالت عبارت است از تأثیر داشتن یک متغیر تبیینی بر روی متغیر واپسیه در حضور دیگر متغیرهای تبیینی.

مجموعه فرض‌ها برای آمین پارامتر ($P \leq i \leq 1$) عبارتند از:

$$H_{i0}: [\beta_1, \dots, \beta_{i-1}, \beta_i, \beta_{i+1}, \dots, \beta_p] = [\beta_1, \dots, \beta_{i-1}, 0, \beta_{i+1}, \dots, \beta_p]$$

$$H_{il}: [\beta_1, \dots, \beta_{i-1}, \beta_i^*, \beta_{i+1}, \dots, \beta_p] = [\beta_1, \dots, \beta_{i-1}, \beta_i^*, \beta_{i+1}, \dots, \beta_p]$$

که در آن β_i^* مقداری معلوم و تعیین شده است. باید ببینیم که حجم نمونه در این حالت چه خواهد شد. برای روشن تر شدن حالت بالا فرض کنید که می‌خواهید تأثیر میزان دوز مصرفی یک داروی خاص (X_1) را بر پرفشار خونی (Y) نشان دهید و از طرف دیگر، نشان دادن رابطه بین پرفشار خونی با تری‌گلیسیرید (X_2) و کلسترول آن (X_3) نیز مورد نظر است. از مطالعات قبلی تیز چنین نتیجه‌گیری شده است که رابطه معناداری بین میزان دوز مصرفی این دارو و میزان تری‌گلیسیرید و کلسترول خون وجود دارد. در اینجا با مسئله‌ای مشابه با آنچه در آزمون فرض‌های بالا ترتیب داده شد روبرو هستیم؛ یعنی:

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

به دلیل همبستگی بین X ها، حجم نمونه‌ای که از رابطه ۶ براساس فقط یک متغیر (مثلاً دوز مصرفی) به دست می‌آید و اگر با این تعداد مدل، برآش شود همه ضرایب رگرسیونی در سطح معناداری α لزوماً توان $1-\beta$ را به دست خواهد داد [۳].

بنابراین باید تدبیلی را بر روی آن انجام داد تا اثر سایر متغیرها نیز در آن ظاهر شود.

بدون این‌که از کلیت مسئله کاسته شود β_1 را در نظر می‌گیریم و فرض می‌کنیم که β_1 براورد کمترین مربعات آن باشد. وایت مور (Whittemore) نشان داده که برای X ها پیوسته و نرمال می‌توان واریانس b_1 ($Var_1(\beta_1)$) را از مدل رگرسیونی تک متغیره ۲ به دست آورد و سپس با استفاده از عامل تورم واریانس (Variance Infusion Factor=VIF)، که در رگرسیون چندگانه قابل محاسبه است (پیوست ۱) واریانس β_1 را در مجموعه چند متغیره P مشکل از متغیر تبیینی (b_1) به دست آورد. پس از آن، حجم نمونه

$$VIF = \frac{Var_p(b_{\backslash})}{Var_{\backslash}(b_{\backslash})} \approx \frac{1}{1-R^2_{X_{\backslash}|X_{\backslash}, \dots, X_p}} \quad (A)$$

ضریب تعیینی است که از رگرسیون $R_{x_1, x_2, \dots, x_p}^2$ بر مجموعه X_1, X_2, \dots, X_p به دست می‌آید. یعنی معادله‌ای به صورت $X_1 = \alpha_0 + \alpha_1 X_2 + \dots + \alpha_p X_p + c'$ داریم و معادله رگرسیونی تک متغیره به صورت ۲ است. برای سادگی فرض می‌کنیم $\mu_{x_1} = 0$ (در غیراین صورت با تبدیلات مناسب می‌توان این کار را انجام داد [۴]). واریانس برآورد کم ترین مربعات پارامتر β_1 (یعنی b_1) از راسته زیر به دست می‌آید [۵]:

$$Var_x(b_i) = \frac{\sigma_x^2}{\sum x_i} \quad (9)$$

حال فرض کنید که متغیر دیگر، یعنی X_2 با $\mu_{X_2} = 0$ نیز به مدل اضافه می‌شود. برآورد ماتریس واریانس کوواریانس پارامترهای برآورده شده به صورت زیر است:

$$Var(b_1, b_2) = \sigma^2 \begin{pmatrix} X^T X \end{pmatrix}^{-1}$$

$$= \sigma^2 \begin{bmatrix} \Sigma X_1^2 & \Sigma X_1 X_2 \\ \Sigma X_2 X_1 & \Sigma X_2^2 \end{bmatrix}^{-1} \quad (1)$$

که در آن X ماتریس متغیرهای تبیینی است ($[X_1 \ X_2]$). از عناصر ماتریس 2×2 $\text{Var}(b_1, b_2)$ می‌توان واریانس b_1 را به صورت زیر به دست آورد [۱]:

$$Var_{\gamma}(b_1) = \frac{\sigma_{\gamma}^2 \Sigma X_{\gamma}^2}{\Sigma X_{\gamma}^2 \Sigma X_{\gamma}^2 - (\Sigma X_{\gamma} X_{\gamma})^2} \\ = \frac{\left(\frac{\sigma_{\gamma}^2}{\sigma_1^2} Var_1(b_1) \right)}{1 - R_{X_1 | X_{\gamma}}^2} \quad (11)$$

مقدار $\left(\frac{\sigma_2}{\sigma_1} \right)$ در بیش تر مواقع کوچک تر یا نزدیک به

یک است و از آن جا که با اضافه شدن یک متغیر تبیینی به مدل، یک درجه آزادی از درجه آزادی مربوط به عبارت مجموع مربعات خطاكم می شود، ممکن است سرآورده

بحث و نتیجه‌گیری

مطالبی که به بحث از آن‌ها پرداختیم از جمله مسائلی هستند که به وفور در تحقیقات گوناگون با آن‌ها مواجه می‌شویم. به دلیل اهمیتی که مدل‌های رگرسیونی در نشان دادن روابط معنادار بین متغیرها دارند و همچنین بنابر تفاوت فاحشی که بین دو حجم نمونه‌ای که به طور مثال ازائه شد وجود دارد، برای دستیابی به نتایج توانمند آماری باید بیش از پیش به ارتباطات خطی بین متغیرهای مستقل در تعیین حجم نمونه توجه گردد؛ زیرا در کلیه مطالعات آماری در سطح معناداری ثابت «با کاهش پیدا کردن حجم نمونه از توان آزمون‌ها کاسته خواهد شد» یعنی احتمال خطای نوع دوم افزایش می‌یابد. همچنین با اعمال تغییر پاسخ گستته نیز سر و کار داریم تعیین داد. البته در این خصوص به دلیل نبودن تعریف مشخص و استاندارد برای ضریب تعیین² با انواع مختلفی از R^2 ها مواجه هستیم که کمی مشکل‌ساز است؛ ولی به هر حال باز هم می‌توان عامل تورم واریانس را با تعاریف مختلف R^2 محاسبه و توان آزمون‌ها را براساس این² R^2 ها با استفاده از شبیه‌سازی‌های کامپیوتری برآورد کرد.

پیوست ۱

روش محاسبه VIF

۱- حالتی که مدل از حالت یک متغیر X_1 به حالت p متغیره X_1, X_2, \dots, X_p افزایش می‌یابد.

فرض کنیم متغیر $Var_p(b_1)$ و $Var_1(b_1)$ به ترتیب واریانس برآورده پارامتر β_1 و β_p باشند که از مدل‌های رگرسیونی با یک و P متغیر تبیینی به دست آمده‌اند. در این صورت با اضافه شدن P-1 متغیر X_2, X_3, \dots, X_p به مدل فقط شامل متغیر X_1 ، واریانس پارامتر برآورد شده $(b_1)_{Var_1}$ تبدیل $Var_p(b_1)$ شده است. به عبارت دیگر، $(b_1)_{Var_1}$ با

$$\text{Var}_x(b_t) = \text{Var}_1(b_1) \times \text{VIF} \quad (N)$$

متورم شده است. این رابطه از همخطی بین متغیرها تیجه‌سنجی گردد [۴ و ۵]. نشان می‌دهیم که در پیش‌تر موقایع که در آن

۲- حالتی که مدل q متغیره به حالت p متغیره ($p > q$) افزایش می‌یابد.

در صورتی که نظیر حالت قبل عمل کنیم، واریانس برآورد (b_1) در مدل q متغیره $\text{Var}_q(b_1)$ و در مدل p متغیره توسط رابطه زیر با یکدیگر مرتبط خواهد بود:

$$\begin{aligned} \frac{\text{Var}_p(b_1)}{\text{Var}_q(b_1)} &= \frac{\left(\frac{\sigma_p^2}{\sigma_q^2}\right) \left(1 - R_{X_1 | X_2, \dots, X_p}^2\right)}{1 - R_{X_1 | X_2, \dots, X_q}^2} \\ &\leq \frac{1 - R_{X_1 | X_2, \dots, X_p}^2}{1 - R_{X_1 | X_2, \dots, X_q}^2} \\ &= \frac{1}{R_{(X_1, X_{q+1}, \dots, X_p) | (X_2, X_3, \dots, X_q)}} \end{aligned} \quad (17)$$

که در آن، ضریب تعیین چندگانه (Multiple determination coefficient)

$R_{(X_1, X_{q+1}, \dots, X_p) | (X_2, X_3, \dots, X_q)}$ میزان ارتباط خطی بین متغیرهای تبیینی (X_1, X_{q+1}, \dots, X_p) و (X_2, \dots, X_q) را که از رگرسیون چندمتغیره (X_p, \dots, X_{q+1}, X_1) بر روی (X_2, \dots, X_q) به دست می‌آید، تعیین می‌کند.

در اینجا $\left(\frac{\sigma_p^2}{\sigma_q^2}\right)$ بیش تراز $\left(\frac{\sigma_p^2}{\sigma_1^2}\right)$ به مقدار یک نزدیک است.

- Whittemore, A.: Sample size for Logistic Regression with small Response probability, Journal of the American Statistical Association, 1981; 76, 27-32.
- Sokal, R.R. & Rohlf, F. J. Biometry, Freeman & Company, New York, 1989; 56.
- Liu, Gand, Liyan. K.Y., Sample size Calculation for studies with Correlated Observations, Biometrics. 1997; 53, 537-547.
- Neter, John. & Authors, Applied Linear Statistical Models, IRwin, USA. 1988.

واریانس گاهی اوقات نیز به طور مختصر از یک $\left(\frac{\sigma_p^2}{\sigma_1^2}\right)$ بیش تر شود.

لذا با تعمیم رابطه ۱۱ به P متغیر تبیینی (به جای ۲ متغیر) رابطه ۸ به دست خواهد آمد؛ یعنی:

$$\text{Var}_p(b_1, \dots, b_p) = \sigma_p^2 (X^T X)^{-1} = \sigma_p^2 \Sigma \quad (12)$$

که در آن Σ ماتریس واریانس - کوواریانس متغیرهای مستقل است.

با استفاده از این رابطه، اندرسون (Anderson) نشان داده است [۶] که اگر:

$$\Sigma_{11}^{-1} = \Sigma X_1^T (1 - R_{X_1 | X_2, \dots, X_p}^2) \quad (13)$$

واز رابطه ۱۲ به دست می‌آید:

$$\begin{aligned} \text{Var}_p(b_1) &= \sigma_p^2 \Sigma_{11} = \frac{\sigma_p^2}{\Sigma_{11}} \\ &= \frac{\sigma_p^2}{\Sigma X_1^T (1 - R_{X_1 | X_2, \dots, X_p}^2)} \\ &= \frac{\left(\frac{\sigma_p^2}{\sigma_1^2}\right) \text{Var}_1(b_1)}{1 - R_{X_1 | X_2, \dots, X_p}^2} \end{aligned} \quad (14)$$

$$\text{Var}_p(b_1) \leq \frac{\text{Var}_1(b_1)}{1 - R_{X_1 | X_2, \dots, X_p}^2} \quad (15)$$

در نتیجه حد بالایی برای VIF عبارت است از:

$$\frac{1}{1 - R_{X_1 | X_2, \dots, X_p}^2} \quad (16)$$

وضعیت‌های بسیار نادری نیز پیش می‌آیند که در آنها $\frac{\sigma_p^2}{\sigma_1^2} > 1$ ؛ اما باز هم رابطه ۱۶ تقریب خوبی برای VIF است [۷].

- Statistical Analysis, Wiley, New York, 32, 1958.
7. Guenther, W. C. Sample size formulas for normal theory t tests american statistician, 1981; 35, 243-244.
5. Weisberg, Sanford. Applied Linear Regression, Wiley, New York, 1988.
6. Anderson, T.W. An Introduction to Multivariate

