

تعیین حجم نمونه در مطالعات فاقد تخصیص تصادفی بقاء با داده‌های بدون سانسور و دارای سانسور

دکتر سقوط فقیه زاده^۱، مهدی رهگذر

در مدل‌های نیمه پارامتری متناسب کاکس که واریانس پارامتر به حالت ساده‌تری تبدیل نمی‌شود باید از یک مدل شبیه سازی استفاده شود، وقتی ضریب تعیین رگرسیون نسبتاً بزرگ باشد توانی که براساس آزمون لگ - رتبه‌ای تطبیق نشده محاسبه می‌شود، توان راییش از مقدار واقعی برآوردمی کند ولی وقتی ضریب تعیین نزدیک به صفر باشد تفاوت در توانها کمتر می‌شود. بالغایش ضریب تعیین رگرسیون، تفاوت بین توانهای حاصل از آزمون لگ - رتبه‌ای و ضریب تعیین تطبیق شده باروش مورد نظر افزایش می‌یابد.

- واژه‌های کلیدی: آنالیز بقاء، تعیین حجم نمونه، عدم تخصیص تصادفی، داده‌های فاقد سانسور، داده‌های دارای سانسور.

مقدمه

در مطالعات آماری و تحلیل‌های بقاء با توجه به اهداف مورد بررسی لازم است که نمونه‌ای را بحاجت لازم انتخاب نماییم تا براساس آن بتوان در سطح اشتباہ معین و توان مناسب استنباط‌هایی را در مورد پارامترهای مورد نظر به عمل آوریم. تخصیص تصادفی افراد به تیمارها، لازمه بسیاری از روش‌ها و تحلیل‌های آماری از جمله در آنالیز بقاء است. امام‌طالعاتی نیز وجود دارند که تخصیص تصادفی امکان پذیرنامی باشد. فرض کنید در یک مرکز آموزشی بخواهیم زمان تأوقیع سرطان ریه را، در دو گروه سیگاری و غیر سیگاری‌ها مقایسه کنیم، در اینجا امکان تخصیص تصادفی افراد به دو گروه فوق امکان پذیرنیست یاد مماییه روند آخرین باروری زنان شهری با زنان روستائی، تخصیص زنان به دو گروه شهری و روستائی به صورت تصادفی صورت نمی‌پذیرد. در این مثال ممکن است زمان تأوقیع سرطان ریه به عوامل دیگری غیر از کشیدن سیگار نیز بستگی داشته باشد. بنابراین در رگرسیون بقاء رابطه متغیر مورد نظر Z که تخصیص به گروهها را تعیین می‌کند و دارای اثر ثابت است، با متغیرهای کمکی (Covariates) W و متغیر تصادفی زمان T نه تنها در مطالعه بلکه در برآورد حجم نمونه نیز باید مدنظر قرار گیرد. جهت برآورد حجم نمونه برای مطالعات تصادفی در آنالیز بقاء روش‌های پیشنهاد شده است از جمله توسط فریدمن (Freedman) (۱۹۸۲)، شوئن فلد (Schoenfeld) (۱۹۸۳)، لاجین و فولکس (Lachin and Lan) (۱۹۹۲)، (Foulkes) (۱۹۸۶)، (Lagatos and Lan) (۱۹۹۲) و برای مطالعات فاقد تخصیص تصادفی آنالیز بقاء، فرمولی برای برآورد

چکیده مقاله
مقدمه. در پژوهش‌هایی که از روش‌های آنالیز بقاء بهره گرفته می‌شود، تعیین حجم نمونه کافی برای دستیابی به توان آماری مناسب دارای اهمیت است. در روش‌های پارامتری و ناپارامتری آمار کلاسیک تخصیص تصادفی نمونه ها از شرایط اساسی محسوب می‌شود. تحقق چنین شرطی در مسیاری از پژوهش‌های کارآزمایی بالینی و بهداشتی امکان پذیرنیست از آنجاکه مطالعات مشاهده‌ای بقاء نمی‌توان افراد را بطور تصادفی به تیمارها متسبب کرد لذا می‌توان از رگرسیون خطی شامل متغیر مستقل با اثربخشی در تحلیل استفاده نمود، نوشتار فوق حاصل تعیین حجم نمونه در اینگونه مطالعات آنالیز بقاء با داده‌های فاقد سانسور و دارای سانسور می‌باشد.

روشها. در مطالعات فاقد تخصیص تصادفی بقاء، می‌توان از رگرسیون خطی که شامل متغیر دارای اثربخشی است، استفاده نمود. چنین رگرسیونی در واقع امید ریاضی متغیر وابسته به شرط متغیر مستقل را بیان می‌کند. با در نظر گرفتن متغیر دوچشمی برای تعیین عضویت فرد در یکی از دو گروه مورد مقایسه و منظور کردن توزیع نمایی برای تابع مخاطره در رگرسیون بقاء، تابع درستنمایی تشکیل گردید. با توجه با اینکه توزیع مجذبی برآورده گریز از جمله درستنمایی نرمال می‌باشد؛ فرمول واریانس برآورده گریز ضریب رگرسیون بدست آمد. سپس با اینکه توزیع مجذبی برآورده گریز ضریب رگرسیون بدست آمد، سپس با اینکه توزیع مجذبی برآورده گریز ضریب رگرسیون بدست آمد. سپس تعیین رگرسیون، فرمولهایی برای تعیین حجم نمونه هم در داده‌های فاقد سانسور و هم دارای سانسور بقاء فراهم شده است.

نتایج. برای مقایسه تابع مخاطره مربوط به دو گروه که فاقد سانسور در داده‌ها می‌باشند، برآورده از ضریب تعیین رگرسیون بقاء، نسبت مخاطره و نسبت عضویت در گروه و واریانس‌های آنهاز تابع درستنمایی بدست آمد. در صورتیکه داده‌ها دارای سانسور باشند علاوه بر عناصر فوق برآورده از احتمال سانسور باید مدنظر قرار گیرد. با اینکه واریانس برآورده گرماکریم درستنمایی و با توجه به توزیع مجذبی نرمال آن و به کمک ضریب تعیین رگرسیون فرمولهایی برای تعیین حجم نمونه بدست آمده است. حجم نمونه تعیین شده توان لازم جهت آزمونی که برای اثر سایر متغیرهای مستقل تطبیق یافته است را دربردارد.

بحث. کاربرد مدل رگرسیونی در مطالعات فاقد تخصیص تصادفی بقاء موجب شده است که فرمولهای مناسبی برای برآورده حجم نمونه، هم در مطالعات فاقد سانسور و هم در مطالعات دارای سانسور فراهم شود که در سطح اشتباہ α دسترسی به توان آماری لازم نیز ممکن می‌گردد.

۱- دانشیار گروه آمارزیستی، دانشگاه تربیت مدرس

می‌باشد. آزمونی باسطح اشتباه α به ازای مقادیر زیر، فرض H_0 را رد می‌کند:

$$|\hat{\beta}_1| > Z_{1-\alpha/2} \sqrt{var(\hat{\beta}_1)}$$

که در آن $Z_{1-\alpha}$ چندک $(1-\alpha)\%$ توزیع نرمال استاندارد است. برای داشتن توان $\gamma = 1 - \alpha$ تحت فرض

$$Pr\left\{ |\beta_1| > Z_{1-\alpha/2} \sqrt{var(\hat{\beta})} \right\} = 1 - \gamma, \quad H_a: \beta_1 \neq 0$$

يعني:

واريانس مجاني $\hat{\beta}_1$ می‌شود:

$$var(\hat{\beta}_1) = \left[\frac{\beta_1}{Z_{1-\alpha/2} + Z_{1-\gamma}} \right]^2 \quad (1)$$

که در آن $Z_{1-\alpha}$ چندک $(1-\alpha)\%$ توزیع نرمال استاندارد می‌باشد. توجه شود که $var(\hat{\beta}_1)$ بستگی به تعداد افراد در مطالعه یعنی n دارد ولی اين وابستگی در معادله (1) محوشده است. هم برای داده‌های دارای سانسور و هم در داده‌های فاقد سانسور نشان داده می‌شود که $var(\hat{\beta}_1)$ می‌تواند به صورت تابعی از ضریب تعیین رگرسیون خطی $R^2_{z^*w^*}$ بحسب W بیان شود. بامسالوی قراردادن واريانس حاصل از رگرسیون با معادله (1) فرمولی برای n ، یعنی حجم نمونه مورد لزوم برای آزمونی در سطح اشتباه α با توان $\gamma = 1 - \alpha$ تیجه خواهد شد.

داده‌های فاقد سانسور در داده‌های فاقد سانسور در داده‌های فاقد سانسور $E(X_i) = \frac{1}{\lambda_i}$ است پس از یافتن $var(\hat{\beta}_1)$ به صورت زیر:

$$var(\hat{\beta}_1) = \left\{ \left(1 - R_{z^*w^*} \right) n P_1 (1 - P_1) \right\}^{-1}$$

باتوجه به توزیع مجاني نرمال (β_1) تعداد نمونه لازم برای انجام آزمون در سطح اشتباه α با توان $\gamma = 1 - \alpha$ خواهد شد:

$$n = \left[\frac{Z_{1-\alpha/2} + Z_{1-\gamma}}{\beta_1} \right]^2 \frac{1}{P_1 (1 - P_1) \left\{ 1 - R_{z^*w^*} \right\}} \quad (2)$$

که در آن P_1 نسبت افراد در گروه 1 می‌باشد. با تغییر ازالت آزمون فرض به سمت برآورد فاصله‌ای، حجم نمونه لازم برای برآورد β_1 توسيع يك فاصله اطمینان با عرض $2d$ می‌شود:

$$var(\hat{\beta}_1) = \left[\frac{\beta_1}{Z_{1-\alpha/2} + Z_{1-\gamma}} \right]^2 \quad (3)$$

که در آن d مقدار ثابتی است که از پيش معين می‌شود.

داده‌های دارای سانسور در داده‌های دارای سانسور $E(X_i) = \frac{1}{\lambda_i} (1 - \pi_i)$ می‌باشد، بنابراین، واريانس خواهد شد:

$$var(\hat{\beta}_1) = \left\{ (1 - \pi_i) \left(1 - R_{z^*w^*} \right) n P_1 (1 - P_1) \right\}^{-1}$$

حجم نمونه توسط برناردو و همکاران (M.V.patricia Bernardo) (۱۹۹۵) (۲۰۰۰) ارائه شده (۵) و در آن فرض شده است که توزیع زمان تا وقوع حادثه یک توزیع نمایی است. آنها به تبعیت از روش رگرسیونی که توسط لیپزیتس و پارزن (Lipsitz and parzen) (۱۹۹۵) برای مطالعات غیر مشاهده‌ای ارائه شده از رگرسیون نمایی جهت مدل سازی مخاطره Z به W استفاده نموده‌اند. عمولاً در مقایسه‌های دارای تخصیص تصادفی بقاء توان آزمون براساس فرمول آزمون لگ-رتبه‌ای (Log-rank) محاسبه می‌شود. برناردو و همکاران (۲۰۰۰) (۵) نشان داده‌اند که توان بدست آمده براساس فرمول آزمون لگ-رتبه‌ای، بيش از مقدار واقعی آن برآورد می‌شود. وايت مر (Whittemore) (۱۹۸۱) (۷) مساله مشابه‌ای را برای توزیع دوچمله‌ای و سیگنورینی (Signorini) (۱۹۹۱) (۸) برای داده‌های شمارشی موردمطالعه قرار داده‌اند. در این نوشتار بالاستفاده از روش برناردو و همکاران (۲۰۰۰) (۵) ولاپیزیتس و پارزن (۱۹۹۵) (۶) فرمول‌های تعیین حجم نمونه در مطالعات آنالیز بقاء با داده‌های فاقد سانسور و دارای سانسور و کاربردی از آنها ارائه می‌شود.

روش‌ها

لیپزیتس و پارزن (۱۹۹۵) محاسبه حجم نمونه را در مطالعات غیر تصادفی مورد بحث قرار داده‌اند که در مقایسه دوتیمار وقتی افراد به صورت تصادفی به دوتیمار تخصیص داده نمی‌شوند، می‌توان از رگرسیون چندگانه متغیر وابسته به شرط متغیر کمکی استفاده نمود (۶). برناردو و همکاران (۲۰۰۰) این روش را به آنالیز رگرسیون بقاء، تعیین داده فرمولهایی را اثبات نموده‌اند (۵).

زمینه موضوع در آنالیز بقاء

فرض کنید $i=1, \dots, n, T$ یک متغیر تصادفی نامتفق باشد که برای n فرد مورد نظر، زمان را لحظه ورود در مطالعه تا وقوع حادثه مورد نظر بیان می‌نماید و C_i متغیر تصادفی میین سانسور شدن باشد که توزیع آن دلخواه است و X_i به صورت (C_i, T_i) \min یعنی Z_i یک متغیر کمکی دوچمله‌ای است که عضویت در گروه 1 در یک مطالعه غیر تصادفی مشخص می‌کند. اگر فرد i در گروه 1 باشد مقدار آن به صورت $=1$ و اگر فرد i در گروه 2 قرار گیرد، $=0$ است. سایر متغیرهای کمکی باتام W_i برای هر فرد به صورت بردار منظور و درایه (base-line) قرار می‌گیرند.

فرض می‌کنیم که T_i به صورت شرطی مستقل از C_i به شرط Z_i و W_i باشد. هدف اصلی مقایسه تابع مخاطره مربوط دو گروه است، در رگرسیون نمایی، تابع مخاطره فرد i به صورت زیراست:

$$\exp(\hat{\beta}_0 + \hat{\beta}_1 Z_i + \hat{\beta}_2 W_i)$$

بنابراین توان آزمونی به صورت $H_0: \beta_1 = 0$ بستگی به واريانس $\hat{\beta}_1$ خواهد داشت. با استفاده از نظریه درستنمایی استاندارد و برای مقادیر به اندازه کافی $var(\hat{\beta}_1)$ بزرگ $\hat{\beta}_1$ دارای توزیع تقریباً نرمال با میانگین β_1 و واريانس (β_1)

در صورتیکه نسبت مردان در جامعه دارندگان سلطان گلوبولهای سفید $P_1 = 0.65$ باشد باتوجه به اینکه $Z_{1-\alpha/2} = 1.96$ و $Z = 1.7$ است باستفاده از فرمول (۲) حجم نمونه برابر ۹۹ نفر خواهد شد.

ب) داده‌های دارای سانسور

پژوهشگری می‌خواهد تابع مخاطره دریافت کنندگان کلیه را دریک مرکز تحقیقاتی کلیه در میان دوگروه شهری و روستایی مقایسه نماید. اگر توزیع بقاء نمایی باشد و احتمال سانسور براساس تجربه‌های قبلی $\pi = 0.3$ و مقدار ضریب تعیین $R^2_{z/w} = 0.92$ باشد، حجم نمونه لازم جهت کسب توان ۹۰ درصد در سطح اشتباه $\alpha = 0.05$ چقدر است؟ در صورتی که نسبت شهری‌ها به روستایان در دریافت کنندگان کلیه در مراکز برابر $0.75 = p_1$ باشد، با توجه به این که $Z_{1-\alpha/2} = 1.96$ و $Z_{1-\beta} = 0.87$ است با استفاده از فرمول (۴) حجم نمونه برابر ۱۸۸ نفر خواهد شد. در داده‌های موجود با منظور کردن نسبت مخاطره $\beta_1 = 0.22$ در رگرسیون مخاطره متناسب یانمایی، توان‌ها محسوب شدند، جدول ۱ توان حاصل از روش آزمون تطبیق نشده (Unadjusted test) دو نمونه‌ای لگ - رتبه‌ای و روش آزمون والد (Wald) در رگرسیون تطبیق شده رابه ازای نسبت مخاطره فوق نشان می‌دهد. ضرایب تعیین با انجام $R^2_{z/w}$ رگرسیون هر متغیر بر حسب سایر متغیرهای موجود بدست آمده است و P_1 عبارت از نسبت تعلق افراد به گروه اول می‌باشد.

لذا حجم نمونه برای آزمونی در سطح اشتباه α با توان $\gamma = 1$ در این حالت می‌شود:

$$(4) \quad n = \frac{1}{\left[\frac{Z_{1-\alpha/2} + Z_{1-\gamma}}{\beta_1} \right]^2 \frac{1}{p_1(1-p_1) \left[1 - R^2_{z/w} \right] (1-\pi)}}$$

که π عبارت است از احتمال اینکه یک فرد به طور تصادفی ضمن نادیده انگاشتن متغیرهای کمکی دچار سانسور شود.

یافته‌ها

حال با استفاده از داده‌های مطالعاتی قبلی بقاء (۹). کاربرد فرمولهای تعیین حجم نمونه مذکور ارائه می‌شود.

الف) داده‌های فاقد سانسور

محققی می‌خواهد تابع مخاطره تسکین استروئید - ایندیوسد (Steroid-induced) راکه در اثر استفاده از داروی مرکاپتوپورین - ۶ (Mercaptoperine) پدیده می‌آید در بین زنان و مردانی که دچار سلطان گلوبولهای سفید هستند مقایسه نماید. در صورتیکه بداند سانسور داده‌ها اتفاق نخواهد افتاد و توزیع داده‌ها نمایی است و مطالعات قبلی مقدار ضریب تعیین را $R^2_{z/w} = 0.85$ و $\beta_1 = 0.45$ نشان داده باشد. حجم نمونه لازم جهت کسب توان ۸۵ درصد در سطح اشتباه $\alpha = 0.05$ چقدر است؟

جدول ۱: محاسبه توان براساس آزمون دونمونه‌ای تطبیق نشده لگ - رتبه‌ای و روش تطبیق شده برای آزمون دوطرفه با خطای 0.05

نوع داده	متغیر مورد نظر	مقدار	توان		
فاقد سانسور	سن	۰/۰۱۹	۰/۴۴	۰/۲۲	۸۱/۲۲
	نگاریتم تعداد گلوبولهای سفید	۰/۰۶۴	۰/۳۵	۰/۲۲	۷۵/۷۱
دارای	وضعیت انجام عمل	۰/۲۰	۰/۲۹	۰/۲۲	۶۹/۴۸
سانسور	مرحله بیماری	۰/۱۸۵	۰/۷۸	۰/۲۲	۸۰/۹۰
	داشت علائم سیستمیک	۰/۲۲	۰/۴۹	۰/۲۲	۸۴/۰۹

تصادفی، در تحلیل رگرسیون خطی، متغیر مستقل نقش عدم تخصیص تصادفی را ایفاء می‌نماید. این مطلب به آنالیز بقاء نیز تعیین داده می‌شود. به طوریکه اگر بخواهیم تابع مخاطره دوگروه راکه امکان تخصیص تصادفی برای آنها وجود ندارد باهم مقایسه نمائیم کافی است یک متغیر Z را طوری در نظر بگیریم که مقادیر آن میان گروهها باشد و سایر متغیرهای تبیینی (مستقل) رابه صورت بردار W در مدل رگرسیون منظور می‌نماییم. استفاده از مدل رگرسیونی در مطالعات غیر تصادفی آنالیز بقاء موجب می‌شود که فرمولهای مناسبی برای برآورد حجم نمونه هم برای مطالعات فاقد سانسور و هم دارای سانسور فراهم شود. در این صورت در سطح اشتباه α پیدا کردن حجم نمونه برای رسیدن به توان آماری لازم نیز ممکن می‌گردد. اگر β_1 ضریب متغیر Z در مدل رگرسیون خطی بیان می‌شود بنابراین برای تحلیل داده‌های فاقد تخصیص

بحث هرچند در روش‌های پارامتری و ناپارامتری در امار کلاسیک انتخاب تصادفی نمونه‌ها در بسیاری از مدل‌های آماری ارزشیابی اساسی محسوب می‌شود. لیکن در بسیاری از پژوهش‌های پژوهشگرانی و بهداشتی امکان تخصیص تصادفی نمونه‌ها فراهم نیست و همواره این سوال مطرح بوده است که در تحلیل چنین داده‌هایی باید از چه روش‌هایی استفاده نمود. بانک اساسی که آنالیز رگرسیون خطی ایفاء می‌نماید اکنون به این نیاز پاسخ مناسبی داده شده است. از آنجا که در تحلیل مدل رگرسیون خطی چند گانه با متغیر مستقل دارای اثر ثابت، متغیر پاسخ (وابسته) و با قیمانده‌ها دارای توزیع تصادفی نرمال بوده و امید ریاضی متغیر وابسته به شرط متغیر مستقل به صورت رگرسیون خطی بیان می‌شود بنابراین برای تحلیل داده‌های فاقد تخصیص

طبق نظریه درستنمایی ، که از حل معادلات امتیاز (score equations) بدست می‌آید. برآوردهای سازگار و مجانبآ نرمال برای β وجود خواهد داشت که دارای ماتریس کوواریانس $E(\beta)\beta E(\beta)^{-1}$ باشد.

داده‌های فاقد سانسور
در داده‌های بدون سانسور $E(X_i) = \frac{1}{n} \sum_{j=1}^n Z_j$ است. پس از بدست آوردن ماتریس مربوط به $E(\beta)$ واریانس مجذوبی بستگی به عنصر (۲۰۲)، در معکوس ماتریس مذکور دارد. لذا با کمک قانون ماتریسهای مجذوبی، وطبق روابط ماتریسها، می‌توان $\text{var}(\beta) = \text{var}(Z) / n$ صورت زیرنوشت (۱۰) را:

$$Z^* \left\{ I - W^* \left[W^{*'} W^* \right]^{-1} W^{*'} \right\} Z^* = Z^* \left(1 - \frac{1}{n} J \right) Z^* - Z^* \left\{ W^* \left[W^{*'} W^* \right]^{-1} W^{*'} - \frac{1}{n} J \right\} Z^*$$

که در آن Z^* برداری $n \times 1$ است که عنصر آن Z_i^* باشد و W^* یک ماتریس $n \times p$ است که سطر آن W_i^* است. و J ماتریس تمام عناصر آن ۱ هستند. با توجه به رابطه اساسی موجود درنظریه رگرسیون است وقتی که رگرسیون Z^* بر حسب W^* بیان شود [۱۱] [۱۱] می‌توان رگرسیون Z^* را با $SS_{\text{reg}} = SS_{\text{tot}} - SS_{\text{res}}$ حداکثری تهیه کرد: (۵)

$$Z^* \left\{ I - W^* \left[W^{*'} W^* \right]^{-1} W^{*'} \right\} Z^* = \left(1 - R_{z^*|w^*}^2 \right) \sum_{i=1}^N (Z_i - \bar{Z})^2$$

که در آن $R_{z^*|w^*}^2$ ضریب تعیین حاصل از رگرسیون Z^* به شرط W^* است. چون Z_i یک متغیر تصادفی دوچمלה‌ای است، معادله (۵) ساده شده به صورت زیر در می‌آید:

$$\hat{\beta} = Z^* \left\{ I - W^* \left[W^{*'} W^* \right]^{-1} W^{*'} \right\} z^* = \left(1 - R_{z^*|w^*}^2 \right) n P_1 (1 - p_1) \quad (6)$$

با مساوی قراردادن معادلات (۱) و (۶) و حل آن برای β تعداد افراد لازم برای انجام آزمونی درسطح اشتباه α با توان $1 - \beta$ به صورت رابطه (۲) در می‌آید.

بافرض اینکه (Z_i^*, W_i^*) مشاهدات مستقل با توزیع یکسان باشند، $R_{z^*|w^*}^2$ در احتمال به مریع ضریب همبستگی چندگانه $\rho_{z^*|w^*}^2$ همگراست (کریستنسن) (Christensen) (۱۹۷۲) (۱۲). بنابراین، در مطالعات آینده نگر که در آن دارای ماتریس طرح (Design matrix) نیستیم در واقع $\rho_{z^*|w^*}^2$ را برابر آورد می‌کنیم و نه مقدار نمونه‌ای آن $R_{z^*|w^*}^2$ را.

داده‌های دارای سانسور
فرض کنید که افراد بعداز سال وارد مطالعه می‌شوند، و تعقیب (Follow-up) آنها تا سال دیگر پس از این اطلاع مطالعه ادامه می‌یابد. فرض کنید $f(c_i)$ می‌باشد از آنجاکه از یک توزیع نمایی پیروی می‌کند. لذا:

$$\pi_i = Pr(T_i > C_i) = \int_0^\infty pr(T_i > C_i | C_i = c_i) f(c_i) dc_i = \int_0^\infty exp(-\lambda_i c_i) f(c_i) dc_i \quad (7)$$

ضریب تعیین $R_{z^*|w^*}^2$ واریانس پارامتر β ، فرمول هایی برای برآورد حجم نمونه به اثبات می‌رسد. فرمول آزمون لگ-رتبه‌ای برای مقایسه تابع مخاطره بین گروههای زمانی مورد استفاده قرار می‌گیرد که تخصیص افراد به گروههای به صورت تصادفی باشد. در صورت فقدان تخصیص تصادفی ممکن است که تابع مخاطره درین گروهها یکسان نباشد لذا روش این مقاله مناسب است. در مدل‌های رگرسیون مخاطره‌های نیمه پارامتری متناسب کاکس که واریانس $\text{var}(\beta)$ به حالت ساده‌تری تبدیل نمی‌شود باید از یک مدل شبیه سازی استفاده شود.

جدول ۱ نشان می‌دهد که وقتی $R_{z^*|w^*}^2$ برای تعیین فرمول حجم نمونه مناسب نسبتاً بزرگ می‌باشد توانی که براساس آزمون لگ-رتبه‌ای تطبیق نشده محاسبه می‌شود، توان را بیش از مقدار واقعی برآورده می‌کند ولی وقتی $R_{z^*|w^*}^2$ نزدیک به صفر باشد مثل متغیرهای سن و لگاریتم تعداد گلولهای سفید، تفاوت در توان کمتر می‌شود. به طور مثال در مقایسه متغیر مرحله بیماری، توان براساس لگ-رتبه‌ای برابر $80/90$ درصد است ولی توان واقعی $72/20$ درصد می‌باشد. براساس جدول ۱، وقتی مقدار $(Z_i - \bar{Z})^2$ افزایش می‌یابد تفاوت بین توانهای حاصل از محاسبات لگ-رتبه‌ای و $R_{z^*|w^*}^2$ تطبیق شده نیز افزایش می‌یابد.

ضمیمه

اکنون اثبات فرمولها که بیشتر براساس مطالعه برناردو و همکارانش (۲۰۰۰) (۵) ولیزیتس و پارزن (۶) استوار است به اختصار بیان می‌شود:

مدل رگرسیون نمایی
اگر تمام متغیرهای کمکی به صورت تصحیح شده با میانگین درنظر گرفته شوند، بعضی از فرمولهای بعدی ساده‌تر می‌شوند. بنابراین مدل رگرسیون نمایی رابه صورت زیر بازنویسی می‌کنیم.

$$\lambda_i = \exp \left[\beta_0 + \beta_1 Z_i + \beta_2 W_i \right]$$

$$\beta^* = \beta_0 + \beta_1 \bar{Z} + \beta_2 \bar{W}, Z_i^* = Z_i - \bar{Z}, W_i^* = W_i - \bar{W}$$

که در آن: $\beta = (\beta_0, \beta_1, \beta_2)$ ، درستنمایی برای β بصورت:

$$L(\beta) = \prod_{i=1}^N \left\{ \lambda_i \exp(-\lambda_i x_i) \right\}^{\delta_i} \exp(-\lambda_i x_i),$$

است که در آن، اگر $T_i \leq C_i$ باشد $\delta_i = 1$ است و در غیر اینصورت $\delta_i = 0$ می‌باشد. نخست ماتریس اطلاع مشاهده شده β از دیربط را ازتابع درستنمایی بدست، می‌آوریم.

ازداده‌های معین سانسور شده‌اند برآورده نمود.

که در آن: $\lambda_i = \exp(\beta_0^* + \beta_1 Z_i^* + \beta_2 W_i^*)$

یادآوری می‌شود که در اینجا:

$$X_i = \min(T_i, C_i)$$

$$E(X_i) = E\left\{ T_i I(T_i < C_i) + C_i I(T_i \geq C_i) \right\}$$

$$= \frac{1}{\lambda_i} (1 - \pi_i)$$

که در آن $I(A)$ یکتابع مشخصه

A است که اگر A

درست باشد مقدار ۱ و در غیر این صورت

مقداره را خیار می‌کند. بنابراین مقدار منتظره

ماتریس اطلاع مشاهده شده نظیر (β)

حاصل از تابع درست‌نمایی است منتها π_i از تمام عناصر ماتریس حذف و تمام

عنصر در مقدار $(1 - \pi_i)$ ضرب می‌شوند.

برآورد احتمال شرطی سانسور شدن i بشرط آنکه متغیرهای کمکی

افراد مشخص باشد از طریق داده‌ها مشکل است. لذا مقدار π_i توسط π

برآورده می‌شود. می‌توان π را به صورت نسبت افرادی که در یک مجموعه

بعد از جایگزینی π توسط ماتریس اطلاع منتظره بدست خواهد آمد.

که وقتی $\pi = 0$ باشد ماتریس حاصل به ماتریس حالت داده‌های فاقد

سانسور تبدیل می‌شود.

بطور مشابه بامساوی قراردادن معادلات، حجم نمونه n برای آزمونی

درسطح اشتباہ α با توان $1 - \gamma$ به صورت رابطه (۴) درخواهد آمد.

توجه نمایید که با انتقال $\pi = 1$ به سمت چپ معادله (۴) فرمولی برای تعداد

شکست‌ها بدست (Failures) خواهد آمد.

برعکس، بامعلوم بودن π ، می‌توانیم توان یک مقایسه را برای آزمونهای

درسطح α به صورت زیر بدست آوریم:

$$1 - \gamma = \Phi \left[\left\{ \beta_0^* np_1 (1 - p_1) \left(1 - R_{z^*}^2 I_w^* \right)^{\frac{1}{2}} (1 - \pi) \right\}^{1/2} - Z_{1-\alpha/2} \right]$$

که در آن Φ تابع توزیع نرمال استاندارد می‌باشد.

در مدل نیمه پارامتری متناسب کاکس مطالعه شبیه سازی که توسط برناردو

و همکاران (۲۰۰۰) (۵) انجام شده نشان داده است که این برآورد

بطور مناسبی عمل می‌کند. بایان مدل رگرسیون کاکس توسط فرایندهای

شمارشی که توسط اندرسون و همکاران (۱۹۸۲) (۱۳) ارائه شده می‌توان به

دست آوردهای جدیدی رسید.

مراجع

- 1- Freedman , L.S.Tables of the number of patients required in clinical trials using the log -rank test. Statist. Med. 1982, 1, 121-129.
- 2- Schoenfeld ,D.A. Sample -size formula for the proportional -hazards regression model . Biometrics . 1983 ,39, 499-503.
- 3- Lachin , J.M.and Foulkes ,M.A. Evaluation of sample size and power for analyses of survival with allowance for nonuniform patient entry , losses to follow -up , noncompliance and stratification . Biometrics. 1986, 42,507-519.
- 4- Lakatos , E.and Lan ,K.K.G.A comparison of sample size methods for the logrank statistic. Statist. Med. 1992, 11, 179-191.
- 5- Bernardo patricia M.V., Lipsitz Stuart R .,Harrington David P.and Catalano paul J.. Sample size calculations for failure time random variables in non-randomized studies. The statistician , 49 , 2000,part 1,31-40.
- 6- Lipsitz, S.R and parzen,M.Sample size calculations for non-randomized studies Statistician 1995, 44,81-90
- 7- Whittemore , A.S. Sample size for logistic regression with small response probability .J.Am Statist .Ass. 1981.76,27-32.
- 8- Signorini, D,F.Sample size for poisson regression. Biometrika.1991,78,446-450.
- 9- Kleinbaum .D.G,. Survival analysis , A self -learning text. Springer-Verlag New York,Inc 1989.
- 10- Snedecor ,G.W.and Cochran,W.G. Statistical Methods. Ames : Iowa State University press 1980.
- 11- Myers , R.H.and Milton , J.S. A 1-1rst coures in the theory of Linear Models .Boston , 1999, PW,S-kent.
- 12- Christensen,R.Plan Answers to complex Questions.New York : Springer .1987.
- 13- Andersen ,P.K. and Gill ,R.D. Cox's regression model for counting processes:a large sample study .Analys of Statist.1982,10,1100-1120.