



دانشگاه فردوسی مشهد
دانشکده علوم ریاضی

پایان نامه جهت اخذ درجه کارشناسی ارشد

آمار ریاضی

طرح آزمایش های بهینه برای مدل رگرسیون پواسون

استاد راهنما

دکتر مجید سرمد

استاد مشاور

دکتر غلامحسین شاهکار

نگارش

منصوره روشن نژاد

شهریور ۱۳۸۸

قدردانی

اینک که در پرتو مهر ایزدی و در جوار بارگاه ملکوتی امام رئوف حضرت رضا (ع) به پایان مرحله‌ای دیگر از دوران تحصیل رسیده‌ام، بر خود لازم می‌دانم از عزیزانی که اینجانب را طی این مسیر یاری نموده‌اند، تشکر و قدردانی نمایم.

در ابتدا از تمامی اساتید ارجمندم در دانشگاه فردوسی مشهد به ویژه جناب آقای دکتر ارقامی که با صبر و دلسوزی مرا راهنمایی نمودند، صمیمانه تشکر می‌نمایم. همچنین مایلم تشکر خالصانه خود را از جناب آقای دکتر سرمد استادیار محترم گروه آمار دانشگاه فردوسی مشهد که راهنمایی این رساله را بر عهده داشتند و مرا در دانش خود سهیم نمودند، ابراز نمایم. همچنین از آقای دکتر شاهکار که سمت مشاور این رساله را بر عهده داشتند تشکر می‌نمایم. همین طور از آقای دکتر طالبی و آقای پورسینا به دلیل وقتی که در اختیارم گذاشتند تشکر می‌نمایم. و مایلم مراتب تشکر و قدردانی خود را از آقایان دکتر هادی جباری نوقابی و دکتر رزمخواه که داوری این رساله را بر عهده داشتند و اینجانب را از راهنمایی‌های خود بهره‌مند ساختند اعلام نمایم. لازم است از منشی محترم گروه آمار سرکار خانم سلیمانی و مسئولین کتابخانه آقای داوودنژاد و سرکار خانم صادقی به دلیل کمک صمیمانه و تهیه مراجع مورد نیاز تشکر ویژه می‌نمایم. همچنین از آقای وطن دوست و خانم نخعی مسئول رایانه دانشکده تشکر و قدردانی می‌نمایم. مایلم از آقای شاهینی و آقای رمضانپور که راهنمایی و مساعدت‌های لازم را به منظور تایپ رساله در محیط تک میکر داشته‌اند و سایر عزیزانی که به نحوی در تهیه این رساله مرا کمک نمودند تشکر نمایم.

از پدر و مادر عزیزم، همسر گرانقدرم و مادر مهربان همسرم که همواره مشوق و حامی من بوده‌اند و با صبر و بردباری بی نظیرشان مرا در طی این مسیر یاری نموده‌اند، خالصانه سپاسگزاری می‌نمایم.

چکیده

در سال های اخیر، توجه به طرح آزمایش های بهینه در مدل های خطی تعمیم یافته افزایش یافته است. با وجود این، بیشتر تحقیقات جاری روی مدل داده های دودویی، مخصوصاً مدل رگرسیون لجستیک یک متغیره مرتبه اول متمرکز است. پژوهش حاضر، این موضوع را روی مدل داده های شمارشی تعمیم داده است. هدف اصلی این تحقیق توسعه و بسط طرح آزمایش های کارا و استوار روی مدل رگرسیون پواسون در مطالعات سم شناسی است.

طرح های D -بهینه برای مدل یک سمی مرتبه اول و مدل دو سمی با اثر متقابل بررسی شده و وابستگی آنها به پارامترهای مدل ارزیابی شده است. کاربرد طرح های D -بهینه، به دلیل این که طرح های بهینه برحسب سطوح دوز مؤثر (ED) به پارامترهای مجهول بستگی دارند، محدود است. بنابراین، برخی از این طرح های عملی مانند طرح های فضای برابر و طرح های D -بهینه شرطی که برحسب سطوح دوز مؤثر بوده و مستقل از پارامترها می باشند، بررسی شده اند. این طرح های عملی وقتی که فضای طرح محدود باشد، کاملاً کارا هستند.

طرح های تهیه شده برحسب سطوح دوز مؤثر شبیه طرح های D -بهینه در برابر عدم اطلاع از پارامترها استوار نیستند. برای بررسی این مسئله، طرح های دنباله ای برای مدل های رگرسیون پواسون پیشنهاد شدند. طرح های کاملاً دنباله ای مورد بررسی قرار گرفتند و به گونه ای تهیه شدند که نسبت به عدم اطلاع از پارامترها، کارا و استوار باشند.

واژه های کلیدی: طرح، رگرسیون پواسون، مدل های خطی تعمیم یافته

پیش‌گفتار

در بسیاری از موارد به خصوص در علوم مهندسی متغیر پاسخ متغیری پیوسته است. معمولاً بر روی این متغیرها مدل‌های خطی به خوبی قابل برازش است. برای چنین مدل‌هایی طرح‌های گوناگونی تاکنون ارائه شده‌اند که از آن جمله می‌توان به طرح‌های کسری اشاره نمود. طرح‌های دیگری نیز در صنعت و علوم داروشناسی به صورت وسیعی استفاده می‌شوند که فرضیات اولیه‌ی طرح‌های معمول برای آن‌ها برقرار نیست. یکی از مهم‌ترین مثال‌ها زمانی است که جواب نهایی یک متغیر شمارشی است. این گونه متغیرها به وسیله مدل‌های خطی قابل بیان نیستند و برای بررسی آن‌ها نیاز به مدل‌های خطی تعمیم‌یافته می‌باشد.

با توجه به اهمیت استفاده از مدل‌های خطی تعمیم‌یافته به خصوص در علوم داروشناسی، در فصل دوم به معرفی مدل‌های خطی تعمیم‌یافته، ساختار آن‌ها و روش به‌دست آوردن برآورد پارامترها و ماتریس اطلاع در این مدل‌ها پرداخته‌ایم. سپس مدل پواسون را که در عمل از کاربرد بیشتری در داده‌های شمارشی برخوردار است، مورد بررسی قرار داده‌ایم.

در فصل سوم طرح‌های D -بهینه برای مدل یک سمی مرتبه اول و مدل دو سمی با اثرمتقابل را بررسی کرده‌ایم و سپس به برخی از طرح‌های عملی برای این دو مدل مانند طرح فضای برابر و طرح D -بهینه شرطی پرداخته‌ایم. در نهایت استواری این طرح‌ها را نسبت به بدمشخص‌سازی پارامترها بررسی کردیم.

در فصل چهارم نظریه‌ی هم‌ارزی کلی را شرح داده‌ایم و از آن برای تحقیق بهینگی طرح‌های بهینه موضعی در فصل سوم بهره گرفته‌ایم. این کار لازم است، زیرا برای

به دست آوردن طرح‌های بهینه از روش‌های عددی استفاده شده است و به‌طور تحلیلی آن‌ها را به دست نیاورده ایم، لذا بهینگی آن‌ها زیر سؤال است.

راه‌حل ارائه شده در فصل سوم جایگزینی پارامترهای مجهول مدل با یک مقدار اولیه‌ی خوب است. این راه‌حل به طرح‌های بهینه موضعی منجر می‌شود. در ابتدای فصل پنجم به بررسی راه‌حل ارائه شده در فصل سوم می‌پردازیم. ایرادهای این روش به صورت عمده عبارتند از:

۱. عدم وجود مقدار اولیه‌ی خوب برای پارامترها

۲. عدم استواری این طرح‌ها در انتخاب‌های ناصحیح پارامترها

با توجه به ایرادهای ذکر شده، در فصل پنجم راه‌حلی برای غلبه بر مشکل وابستگی معیارهای بهینگی مطرح می‌شود. در این فصل ابتدا به معرفی طرح‌های بهینه دنباله‌ای می‌پردازیم. سپس این طرح‌ها را برای مدل پواسون توضیح می‌دهیم. با ارائه الگوریتمی و با روش‌های تکراری، طرح بهینه دنباله‌ای را به دست می‌آوریم. (برنامه این الگوریتم در ضمیمه پایان نامه موجود می‌باشد).

فهرست مطالب

۱	تعاریف و مقدمات	۱
۱	۱.۱ مقدمه	۱
۱	۲.۱ نمادگذاری	۱
۴	۳.۱ پیشینه موضوع و انگیزه تحقیق	۴
۵	۴.۱ مروری بر طرح در مدل‌های خطی	۵
۶	۵.۱ بهینگی طرح	۶
۹	۶.۱ طرح‌های بهینه برای مدل رگرسیون لجستیک	۹
۹	۱.۶.۱ طرح‌های بهینه موضعی	۹
۱۱	۲.۶.۱ طرح‌های دنباله‌ای	۱۱
۱۳	۳.۶.۱ طرح‌های بهینه بیزی	۱۳
۱۴	۴.۶.۱ طرح‌های بهینه کم بیشینه	۱۴
۱۵	۷.۱ طرح‌های بهینه برای مدل رگرسیون پواسون	۱۵
۱۵	۱.۷.۱ مدل‌ها، فرضیات و نمادگذاری	۱۵
۱۸	۲.۷.۱ طرح‌های بهینه: مروری بر کارهای انجام شده	۱۸
۱۹	۸.۱ ساختار پایان نامه	۱۹
۲۱	۲ مدل‌های خطی تعمیم‌یافته	۲۱
۲۱	۱.۲ مقدمه	۲۱
۲۲	۲.۲ خانواده نمایی	۲۲
۲۵	۳.۲ ساختار مدل‌های خطی تعمیم‌یافته	۲۵

۲۶	برآورد پارامترها توسط روش حداکثر درستنمایی	۴.۲
۳۰	مدل پواسون	۵.۲
۳۱	برآورد پارامترهای مدل پواسون به شیوه حداکثر درستنمایی	۱.۵.۲

۳ طرح‌های D - بهینه موضعی

۳۷	مقدمه	۱.۳
۳۸	آزمایشات یک-سمی : مدل مرتبه دوم	۲.۳
۳۹	طرح‌های D - بهینه	۱.۲.۳
۴۰	به دست آوردن طرح‌های D - بهینه*	۲.۲.۳
۴۲	طرح فضاهای برابر (هم‌شانس)	۳.۲.۳
۴۲	به دست آوردن طرح‌های D - بهینه*	۴.۲.۳
۴۳	بدمشخص‌سازی پارامتر	۵.۲.۳
۴۴	محاسبات*	۶.۲.۳
۴۵	آزمایشات دو-سمی : مدل با اثر متقابل	۳.۳
۴۷	طرح‌های D - بهینه	۱.۳.۳
۴۸	محاسبات*	۲.۳.۳
۵۲	طرح‌های D - بهینه شرطی	۳.۳.۳
۵۲	محاسبات*	۴.۳.۳

۴ نظریه هم‌ارزی

۵۵	مقدمه	۱.۴
۵۶	نظریه هم‌ارزی	۲.۴
۵۹	بررسی بهینگی طرح‌های D - بهینه موضعی	۳.۴
۶۰	مدل یک سمی مرتبه دوم	۱.۳.۴
۶۱	مدل دو سمی با اثر متقابل	۲.۳.۴
۶۲	خلاصه	۴.۴

۵ طرح‌های دنباله‌ای

۶۷	طرح‌های دنباله‌ای	۱.۵
----	-------------------	-----

۶۸	۱.۱.۵	الگوریتم
۷۰	۲.۵	مثال عددی
۷۰	۱.۲.۵	محاسبات*
۷۸	۲.۲.۵	محاسبه D - کارایی طرح دنباله‌ای

۱

الف برنامه های کامپیوتری*

		۱.الف	برنامه یافتن نقاط طرح D - بهینه برای مدل ۱۷.۱ و کارایی طرح‌های فضای هم‌شانس
۱		
		۲.الف	برنامه یافتن نقاط طرح D - بهینه برای مدل ۱۹.۱ و D - کارایی طرح‌های D - بهینه شرطی
۳		
۷	۳.الف	برنامه حل مثال عددی طرح دنباله‌ای در فصل پنجم

۱۴

ب واژه نامه

لیست تصاویر

۳۱	نمودار $\lambda = \exp(5 - 2x)$	۱.۲
۶۴	بررسی D -بهینگی طرح‌های جدول (۲.۲.۳) (ناحیه اول)	۱.۴
۶۵	بررسی D -بهینگی طرح‌های جدول (۲.۲.۳) (ناحیه دوم و سوم)	۲.۴
۶۶	..	بررسی D -بهینگی طرح‌های جدول (۴.۳) برای مدل دو سمی	۳.۴
۷۷	مثال طرح دنباله ای	۱.۵

لیست جداول

۲۶	توابع پیوند کانونی برای مدل‌های خطی تعمیم‌یافته	۱.۲
۳۴	تعداد موارد آیدز در استرالیا در دوره‌های سه ماهه از ۱۹۸۴ تا ۱۹۸۸	۲.۲
۴۱	طرح‌های D -بهینه برای مدل (۱۷.۱) ($q_1 = 1, p_1 = p_2 = p_3 = \frac{1}{3}$)	۱.۳
۴۳	D -کارایی طرح‌های فضای هم‌شانس برای مدل (۱۷.۱)	۲.۳
۳.۳	D -کارایی طرح‌های فضای هم‌شانس برای مدل (۱۷.۱) تحت بدمشخص‌سازی	
۴۶	پارامترها	
	طرح‌های D -بهینه برای مدل (۱۹.۱) ، ($q_{11} = q_{21} = 1, p_1 = p_2 = p_3 =$)	۴.۳
۵۰		($p_4 = \frac{1}{4}$)
۵۴	D -کارایی طرح‌های D -بهینه شرطی برای مدل (۱۹.۱)	۵.۳
۷۵	مثال طرح‌های دنباله‌ای	۱.۵
۷۶	ادامه مثال طرح‌های دنباله‌ای	۲.۵
۸۰	D -کارایی طرح‌های دنباله‌ای	۳.۵

فصل ۱

تعاریف و مقدمات

۱.۱ مقدمه

در این فصل به اهمیت موضوع و مرور پیشینه و انگیزه تحقیق حاضر می‌پردازیم. سپس مروری بر طرح‌های بهینه در مدل‌های خطی تعمیم‌یافته خواهیم داشت. مراجع [۲]، [۴۰] و [۵۶] از منابع اصلی این فصل هستند.

۲.۱ نمادگذاری

ابتدا معرفی خلاصه‌ای از نمادهای مورد استفاده را خواهیم داشت. متغیر پاسخ را با Y نشان می‌دهیم و یک مجموعه از متغیرهای توضیحی به وسیله x_1, x_2, \dots, x_p نشان داده می‌شود. معادله پیوند به صورت زیر می‌باشد:

$$g(E(Y)) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p.$$

برای متغیرهای پاسخ Y_1, Y_2, \dots, Y_n ، معادله پیوند به شکل ماتریسی زیر در می‌آید:

$$g(E(\mathbf{Y})) = \mathbf{X}\beta,$$

که در آن

$$\mathbf{Y} = \begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix}$$

بردار متغیرهای پاسخ است.

$$g[E(\mathbf{Y})] = \begin{bmatrix} g[E(Y_1)] \\ \vdots \\ g[E(Y_n)] \end{bmatrix} \quad (1.1)$$

نشان دهنده‌ی بردار توابع $E(Y_i)$ می‌باشد که تابع g برای تمام عناصر بردار یکسان است و

$$\beta = \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_p \end{bmatrix}$$

بردار پارامترها است و X ماتریسی است که عناصر آن نشان دهنده‌ی سطح متغیر توضیحی رسته‌ای و یا نمایانگر مقدار متغیر توضیحی پیوسته است. برای یک متغیر توضیحی x (مانند وزن) مدل شامل جمله‌ی βx است که پارامتر β نشان دهنده‌ی تغییر در متغیر پاسخ به ازای تغییر یک واحدی در متغیر x است. برای متغیر توضیحی رسته‌ای، پارامترهایی وجود دارد که سطوح مختلف یک عامل را نشان می‌دهند.

اگر تعداد پارامترهای مدل p باشد و n مشاهده داشته باشیم، آن‌گاه y بردار تصادفی $n \times 1$ می‌باشد، β بردار $p \times 1$ از پارامترها است و X یک ماتریس $n \times p$ از ثابت‌های معلوم می‌باشد. X اغلب ماتریس طرح و $X\beta$ مؤلفه‌ی خطی مدل نامیده می‌شود.

تعریف ۱.۲.۱. تابعی که $E(Y_i) = \mu_i$ را به مؤلفه‌ی خطی $\mathbf{x}'_i\beta$ مربوط می‌سازد تابع پیوند نامیده می‌شود. مانند

$$g(\mu_i) = \mathbf{x}'_i\beta \quad (2.1)$$

که در آن g تابع پیوند نامیده می‌شود.

مدل‌های خطی زیر را در نظر بگیرید:

$$E(Y_i) = \mu_i = \mathbf{x}'_i\beta \quad ; \quad Y_i \sim N(\mu_i, \sigma^2) \quad (3.1)$$

که در آن Y_1, \dots, Y_n متغیرهای تصادفی مستقل می‌باشد و بردار x'_i نشان‌دهنده‌ی i امین سطر ماتریس طرح X است. این مدل معمولاً به صورت ماتریسی زیر نوشته می‌شود:

$$y = X\beta + e, \quad (۴.۱)$$

که در آن

$$y = \begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix}, \quad X = \begin{bmatrix} x'_1 \\ \vdots \\ x'_n \end{bmatrix}, \quad \beta = \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_p \end{bmatrix}, \quad e = \begin{bmatrix} e_1 \\ \vdots \\ e_n \end{bmatrix}$$

و e_i ها متغیرهای تصادفی مستقل و هم‌توزیع با $e_i \sim N(0, \sigma^2)$ ، $i = 1 \dots n$ می‌باشند. مثالی در این مورد می‌تواند رابطه میان وزن زمان تولد نوزاد و سن حاملگی مادر باشد. اکنون حالتی را در نظر بگیرید که Y_i ها دارای توزیع نرمال نباشند و یا رابطه میان متغیرهای پاسخ و متغیرهای توضیحی به صورت مدل (۳.۱) نباشد. مثلاً بررسی مرگ‌ومیر ماهیان (متغیر پاسخ) در مقابل استفاده‌ی مقدار مشخصی از سم (متغیر توضیحی)، در این مثال متغیر پاسخ پیوسته نیست و یک متغیر دوجمله‌ای است. برای بررسی چنین مسائلی نیاز به مدل‌های دیگری با فرضیاتی به غیر از فرضیات مدل خطی (۳.۱) داریم. در این مثال نیاز به مدلی داریم که توزیع متغیرهای پاسخ متعلق به توزیع دوجمله‌ای باشد و علاوه بر این بنا به تجربه می‌دانیم که **تابعی غیرخطی** $E(Y_i) = \mu_i$ را به $x'_i\beta$ مربوط می‌سازد. از جمله‌ی این مدل‌ها می‌توان به مدل دوجمله‌ای که در آن $g(\mu_i) = \ln \frac{p_i}{1-p_i}$ و مدل پواسون که در آن $g(\mu_i) = \ln(\mu_i)$ می‌باشد، اشاره نمود. یکی از ویژگی‌های مدل‌های خطی تعمیم یافته تابع پیوند غیر خطی است، تعریف دقیق این مدل‌ها را تا مروری بر پیشینه موضوع به تعویق می‌اندازیم.

تعریف ۲.۲.۱. به‌طور کلی، دوز مؤثر یا به اختصار ED ، به عنوان دوز یا مقدار سمی تعریف می‌شود که منجر به یک نسبت مشخص از مرگ و میر یا بقاء می‌شود. به عنوان مثال ED_{70} سم A دوزی از سم A است که منجر به ۷۰٪ مرگ و میر می‌شود. برای آزمایشات یک سمی، این تعریف، صریح و روشن است. برای آزمایشات چند سمی، نیاز به تشخیص اثر ترکیبی تمام سموم و ترکیب هر سم به تنهایی داریم. برای این منظور **دوز مؤثر مدل** یا به اختصار MED را تعریف می‌کنیم و به عنوان تمام ترکیبات ممکن دوزهای تمام سموم که منجر به یک مرگ و میر مشخص یا بقای

مشخص می‌شوند، در نظر می‌گیریم. به عنوان نمونه، در یک آزمایش دو سمی، MED_{30} ، تمام ترکیبات ممکن دوزهای سمی است که مرگ و میر ۷۰٪ را ایجاد می‌کنند. هرچند، دوز مؤثر انفرادی یا به اختصار IED یک سم، دوزی از سم است که منجر به مرگ و میر یا بقای مشخص است. روشن است برای آزمایشات تک سمی MED و IED یکسان هستند.

۳.۱ پیشینه موضوع و انگیزه تحقیق

محققانی را در نظر بگیرید که بر روی داروی ضد سرطان مطالعه و بررسی می‌کند و بنا به تجربیات قبلی می‌داند که تعداد کلونی‌های سرطانی، معمولاً یک متغیر پواسون در نظر گرفته می‌شود. برای استفاده‌ی صحیح از داروی مورد نظر می‌خواهد مدل رگرسیون پواسون را که رابطه‌ی میان تعداد کلونی‌ها و غلظت یا دوز دارو را توصیف می‌کند، مشخص نماید. با توجه به مدل (۳.۱)، واضح است که مسأله در فرضیات مدل‌های خطی صدق نمی‌کند و یکی از مسائل مدل‌های خطی تعمیم‌یافته است. از این نوع مسائل در داروشناسی و سم‌شناسی فراوان است.

مدل‌های خطی تعمیم‌یافته برای داده‌های دودویی و شمارشی، در کتاب‌ها و مقالات بسیاری بررسی شده‌اند [۴۰]. اخیراً، طرح آزمایش‌ها برای مدل داده‌های دودویی، مخصوصاً برای مدل رگرسیون لجستیک، که کاربرد فراوانی در علوم زیست‌شناسی و داروشناسی دارد، توجه زیادی را به خود معطوف ساخته و بسیاری از کتاب‌ها و مقالات، به این مبحث پرداخته‌اند. با وجود این، کار کمی بر روی طرح آزمایش‌های کارا برای مدل داده‌های شمارشی مانند مدل رگرسیون پواسون انجام شده است. در این تحقیق، سعی داریم این شکاف را به وسیله توسعه طرح آزمایش‌های کارا و استوار برای مدل رگرسیون پواسون پرکنیم.

یک بررسی سم‌شناسی شامل اجرای غلظت‌های مختلف سموم بر روی ارگانیزم‌ها و اندازه‌گیری اثر شاخص‌هایی همچون مرگ و میر و رشد و صدمه به تولد و تناسل و... را در بر می‌گیرد. (معمولاً این پاسخ‌ها نمی‌تواند به‌طور مناسبی به عنوان متغیرهای

تصادفی نرمال مدل‌بندی شود). شاخص‌هایی همچون مرگ و میر، پاسخ‌های دودویی و شاخص‌هایی همچون صدمه به تولد و تناسل پاسخ‌هایی شمارا دارند. در چنین موقعیت‌هایی، مدل‌های خطی کلاسیک عملی نیست و معمولاً مدل‌های خطی تعمیم‌یافته برای تحلیل به کارگرفته می‌شود. مثال‌های مشابه فراوانی در کتاب‌ها وجود دارد [۳۵].

۴.۱ مروری بر طرح در مدل‌های خطی

در این بخش به معرفی طرح‌ها در مدل‌های خطی می‌پردازیم تا خواننده بتواند ارتباط مؤثری بین این طرح‌ها و طریقه به‌دست آمدن آن‌ها و تفاوت آن‌ها با مدل‌های خطی تعمیم‌یافته برقرار کند. در این قسمت مدل‌هایی به شکل زیر را در نظر می‌گیریم:

$$y = X\beta + e, \quad (5.1)$$

که در آن

$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}, \quad X = \begin{bmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1p} \\ 1 & x_{21} & x_{22} & \cdots & x_{2p} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \cdots & x_{np} \end{bmatrix}, \quad \beta = \begin{bmatrix} \beta_0 \\ \vdots \\ \beta_p \end{bmatrix}, \quad e = \begin{bmatrix} e_1 \\ \vdots \\ e_n \end{bmatrix},$$

که e_i ها متغیرهای تصادفی مستقل و با توزیع $N(0, \sigma^2)$ $e_i \sim N(0, \sigma^2)$ $i = 1, \dots, n$ می‌باشند. از رگرسیون می‌دانیم که برآوردگر درست‌نمایی ماکسیمم β به صورت زیر به‌دست می‌آید:

$$\hat{\beta} = (X'X)^{-1}X'y. \quad (6.1)$$

از آنجا که برآوردهای به‌دست آمده نااریب هستند، رابطه زیر میان ماتریس واریانس-کوواریانس و ماتریس اطلاع برقرار است:

$$\text{cov}(\hat{\beta}) = \sigma^2(X'X)^{-1} = I^{-1}(X, \beta). \quad (7.1)$$

که در آن I ماتریس اطلاع است.

۵.۱ بهینگی طرح

محققى را که در قسمت قبل توضیح دادیم در نظر بگیرید، او به دنبال کشف مدل صحیح رگرسیون پواسون میان دوز دارو و تعداد کلونی‌های سرطانی بود و برای این هدف بایستی ضرایب مدل را براساس آزمایشاتی برآورد نماید. از آنجا که انجام آزمایشات مستلزم صرف هزینه و زمان است؛ او در نظر دارد برای کشف سریع تر مدل صحیح و صرف هزینه کمتر، آزمایشاتی را انجام دهد که برآوردهای به دست آمده براساس این آزمایشات بسیار به مقادیر واقعی ضرایب مدل نزدیک باشند، همچنین برای آزمون فرضیاتی در مورد ضرایب مدل بهتر است که برآوردها دارای کمترین واریانس باشند و به عبارتی برآوردهای ضرایب رگرسیونى دارای بهینگی در واریانس باشند. بهینگی در واریانس برآوردها که آزمایشگر به دنبال آن است، ویژگی طرح‌های بهینه با معیار D -بهینگی است. برای رسیدن به این هدف ابتدا با طرح‌های بهینه و معیارهای بهینگی طرح آشنا می‌شویم.

نخست معرفی خلاصه ای از معیار بهینگی طرح ارائه می‌کنیم. در سال ۱۹۵۹ کیفر و ولفوویتز [۲۳] یک چارچوب نظری برای معیار بهینگی طرح تهیه کردند به این صورت که طرح را به عنوان اندازه احتمال، بیان کردند، که به هر نقطه‌ی طرح در فضای طرح، مشاهداتی را تخصیص می‌دهد. دسترسى نظری آنها به بهینگی طرح و معرفی آنها از معیار D و E بهینگی برای مدل رگرسیون خطی منجر به پایه‌ای برای معیارهای بهینگی دیگر مانند معیارهای $-Q$ ، $-G$ ، $-F$ و $-A$ -بهینگی گردید. هر معیار بهینگی طرح، نشان‌دهنده هدف ویژه‌ای است که آزمایش به منظور رسیدن به آن اجرا می‌شود یا یک ویژگی مخصوص که در مدل رگرسیون برازش یافته‌ی نهایی وجود دارد.

ایده‌ی اصلی نظریه‌ی بهینگی طرح این است که استنباط آماری درباره‌ی کمیت‌های مورد نظر به وسیله بهینگی انتخاب سطوح متغیرهای کنترل می‌تواند بهبود پیدا کند. به‌طور کلی، معیار بهینگی طرح می‌تواند به صورت معیار برآورد یا معیار پیش‌بینی مشخص شود. یک طرح که از حیث معیار برآورد مانند D -بهینگی، بهینه است، به وسیله مینیمم کردن تغییرات برآوردهای پارامتر، اطلاع پارامتر را ماکسیمم می‌کند.

یک طرح که از نظر معیار پیش‌بینی مانند Q -بهینگی، بهینه است، اطلاع درباره‌ی رویه پاسخ را به وسیله تمرکز روی کمیت پیش‌بینی مدل برازش یافته ماکسیم می‌کند. در سال ۱۹۸۰، سیلوی [۴۷] و در سال ۱۹۹۲، اتکینسون و دنو [۴] اکثر معیارهای بهینگی طرح را که به‌طور معمول استفاده می‌شد، شرح دادند [۵۶]. در این تحقیق، بیشتر روی طرح‌های D -بهینه متمرکز خواهیم شد، لذا D -بهینگی در اینجا با جزئیات معرفی می‌شود.

اگر بردار پارامترهای مدل، به وسیله β نشان داده شود، D -بهینگی اشاره می‌کند که انتخاب طرح باید اطلاع در مورد β را به وسیله مینیم کردن واریانس تعمیم‌یافته برآوردگرهای β ، ماکسیم کند. می‌دانیم که روش‌های معمول برای برآورد β در بسیاری از مدل‌های غیر خطی، برآورد درست‌نمایی ماکسیم (MLE) است. اگر $\hat{\beta}$ برآورد درست‌نمایی ماکسیم β باشد، آنگاه ماتریس واریانس جانبی $\hat{\beta}$ ، متناسب با وارون ماتریس اطلاع فیشر است [۲۷]. این ماتریس اطلاع، پایه‌ای را برای معیار بهینگی طرح به وجود می‌آورد که در (۸.۱) تعریف شده است.

تعریف ۱.۵.۱. اگر شرایط نظم برقرار باشد و بعد β ، $p \times 1$ باشد، ماتریس اطلاع فیشر $I(X, \beta)$ یک ماتریس $p \times p$ است:

$$I(X, \beta) = -E\left[\frac{\partial^2 \log(L(X, \beta))}{\partial \beta \partial \beta}\right], \quad (۸.۱)$$

که در آن $L(X, \beta)$ تابع درست‌نمایی داده‌ها و X ماتریس طرح است. توجه کنید که امید ریاضی تحت Y محاسبه می‌شود.

تعریف ۲.۵.۱. معیار D -بهینگی که واریانس تعمیم‌یافته $\hat{\beta}$ ها را مینیم می‌کند، به‌طور معادل دترمینان ماتریس اطلاع فیشر را ماکسیم می‌کند. معیار D -بهینگی کلی به صورت زیر تعریف می‌شود:

$$\max_{X \in D} \left| \frac{I(X, \beta)}{n} \right|, \quad (۹.۱)$$

که در آن D مجموعه‌ای از تمام طرح‌های ممکن و n اندازه نمونه است. در اکثر مواقع، اندازه نمونه ثابت است، بنابراین، طرح D -بهینه از طریق ماکسیم کردن دترمینان ماتریس اطلاع فیشر، به‌دست می‌آید.

تعریف ۳.۵.۱. از آنجا که طرح D -بهینه، دترمینان ماتریس اطلاع فیشر را ماکسیمم می‌کند، طبیعی است که D -کارایی طرح دلخواه X را چنین تعریف کنیم:

$$Eff_D = \left(\frac{|I(X, \beta)|}{|I(X^*, \beta)|} \right)^{\frac{1}{p}}, \quad (10.1)$$

که در آن X^* طرح D -بهینه با همان اندازه X است و p تعداد پارامترهای مدل است. وقتی D -بهینگی به کار برده می‌شود، D -کارایی تعریف شده در (۱۰.۱) عموماً برای مقایسه‌ی طرح‌های متفاوت استفاده می‌شود. گرفتن نسبت دترمینان‌ها در (۱۰.۱) و به توان $\frac{1}{p}$ رساندن آن‌ها، یک اندازه کارایی را که متناسب با حجم نمونه است و صرف نظر از بعد مدل می‌باشد، نتیجه می‌دهد.

همان‌طور که پیش از این توضیح داده شد، در مدل‌های خطی با واریانس ثابت σ^2 ، ماتریس اطلاع متناسب با ماتریس واریانس - کوواریانس می‌باشد (بنا به رابطه (۷.۱)) برای کاهش واریانس برآوردگرها در مدل‌های خطی ناچاریم که ماتریس $(X'X)^{-1}$ و یا به‌طور معادل I^{-1} را مینیمم نماییم. برای مینیمم کردن وارون ماتریس اطلاع، بایستی دترمینان ماتریس اطلاع را ماکسیمم نماییم، و چون ماتریس اطلاع متناسب با $X'X$ است، کافی است دترمینان $X'X$ را ماکسیمم نماییم. یعنی x ‌هایی را از فضای طرح پیدا نماییم که دترمینان مورد نظر را ماکسیمم نماید. بنابراین در مدل‌های خطی، ماتریس اطلاع به پارامتر مجهول بستگی ندارد و به آسانی می‌توان دترمینان ماتریس اطلاع را براساس x ‌ها ماکسیمم نمود. به عبارت دیگر، ماتریس اطلاع فیشر برای واریانس ثابت مدل خطی به β بستگی ندارد. در نتیجه طرح D -بهینه می‌تواند مستقل از β تعیین و اجرا شود. لکن، برای یک مدل غیر خطی مانند مدل خطی تعمیم‌یافته، ماتریس اطلاع، معمولاً به پارامترهای مدل بستگی دارد. ماتریس اطلاع برای مدل خطی تعمیم‌یافته، به صورت زیر است:

$$I(X, \beta) = X'WX \quad (11.1)$$

که در آن W ماتریس وزن‌دار هسین^۱ است که به پارامترهای مجهول بستگی دارد

^۱Hessian

[۲]. برای یک مدل خطی تعمیم یافته کانونی مانند مدل رگرسیون لجستیک رابطه (۱۱.۱) به صورت زیر تبدیل می شود:

$$I(X, \beta) = X'VX \quad (12.1)$$

که در آن V یک ماتریس قطری است و (i, i) امین عنصر آن، واریانس i امین متغیر است. اما این واریانس ها به β بستگی دارند و لذا ماتریس اطلاع نیز به β وابسته خواهد بود. بنابراین، طرح D -بهینه برای یک مدل خطی تعمیم یافته وابسته به پارامتر است. به منظور یافتن و اجرای طرح D -بهینه، باید پارامترهای مدل، معلوم باشند. این وابستگی به پارامترهاست که مسأله طرح آزمایش ها برای مدل های خطی تعمیم یافته را پیچیده می کند. چندین روش، برای حل این مشکل در کتاب ها پیشنهاد شده که در قسمت بعد مروری بر آنها خواهیم داشت.

۶.۱ طرح های بهینه برای مدل رگرسیون لجستیک

برای حل مسأله آزمایشگری که روی داروی ضد سرطان کار می کند و می خواهد برآوردگرهای ضرایب مدل پواسون را با کمترین واریانس به دست آورد، نیاز به پیدا کردن طرح D -بهینه برای مدل رگرسیون پواسون داریم. همان طور که دیدیم طرح آزمایش های بهینه برای مدل های خطی تعمیم یافته به پارامترهای مجهول بستگی دارد و لذا مسأله ساختاری آنها، لزوماً پیچیده تر از مدل های خطی کلاسیک است. در این بخش، برخی از روش هایی را که برای غلبه بر این مشکل پیشنهاد شده اند، معرفی می کنیم. در کتاب ها و مقالات، بیشتر تحقیق رایج روی طرح آزمایش ها برای مدل های خطی تعمیم یافته، به مدل رگرسیون لجستیک محدود شده است. بنابراین، در توضیح این روش ها، بیشتر در زمینه رگرسیون لجستیک صحبت می کنیم.

۱.۶.۱ طرح های بهینه موضعی

یک شیوهی سادهی مسأله طرح آزمایش ها برای مدل های خطی تعمیم یافته، قبول کردن بهترین حدس برای پارامترها و سپس پیدا کردن نقاطی است که تابع معیار

بهینگی انتخاب شده طرح را در حدس مورد نظر، ماکسیمم کند. این شیوه منجر به طرحی می‌شود که طرح بهینه موضعی^۲ نامیده می‌شود و به وسیله چرنوف [۹] معرفی گردید. بهترین حدس که در طرح بهینه موضعی استفاده می‌شود، ممکن است از آزمایش‌های قبلی بیاید، یا از آزمایشی که مخصوصاً برای این هدف انجام شده، یا صرفاً یک حدس باشد. بدون بحث در نحوه‌ی به‌دست آمدن حدس پارامتر، آن را برآورد اولیه یا حدس اولیه، خواهیم نامید.

اکنون، مروری خلاصه بر طرح‌های بهینه موضعی برای مدل رگرسیون لجستیک ارائه می‌کنیم. بیشتر مطالعات بر روی عنوان مدل رگرسیون لجستیک مرتبه اول یک متغیری متمرکز است که به صورت زیر تعریف می‌شود:

$$y \sim \text{Bernoulli}(P), \quad P = \frac{1}{1 + \exp(-(\beta_0 + \beta_1 x))} \quad (13.1)$$

در مدل (۱۳.۱)، y پاسخ دودویی (۰ یا ۱) است، که با احتمال p مقدار یک را می‌گیرد و p تابعی مستقل از متغیر x است. برای این مدل، کالیش و رزنبرگر [۲۱] یک طرح بهینه دو سطحی را به‌دست آوردند. آن‌ها تقارن حول نقطه $ED_{0.5}$ را تعیین کردند (ED مخفف دوز مؤثر است و ED_{100p} بیانگر مقدار x ای است که با احتمال p تولید می‌شود و دارای پاسخ یک است). طرح D -بهینه دو سطحی به وسیله قرار دادن نیمی از اجراهای آزمایشی در نقطه $ED_{17.6}$ و نیمی دیگر در نقطه $ED_{82.4}$ به‌دست می‌آید. برای تکمیل طرح D -بهینه باید مقادیر $ED_{17.6}$ و $ED_{82.4}$ معلوم باشد، یا به‌طور معادل، باید اطلاعی در مورد β_0 و β_1 داشته باشیم. در عمل، پارامترها ذاتاً مجهول هستند و نیاز به تهیه‌ی چندین حدس اولیه می‌باشد.

کالیش و رزنبرگر [۲۱]، یک طرح دو سطحی G -بهینه را توسعه دادند که ماکسیمم واریانس احتمال پیش‌بینی \hat{P} را روی ناحیه مورد نظر مینیمم می‌کند. ویلیامز [۵۴] در سال ۱۹۸۶ روی معیار بهینگی دیگری که طول بازه اطمینان را براساس برآوردهای درست‌نمایی ماکسیمم می‌کند، کار کرد. کار دیگر روی طرح‌های بهینه موضعی برای مدل رگرسیون لجستیک شامل کار مایرز و هیز [۱۹] است، آن‌ها طرح‌های D -بهینه و Q -بهینه را برای مدل رگرسیون لجستیک دو متغیری بررسی کردند. اگر