



۳۲۷۶۶



دانشگاه صنعتی اصفهان
دانشکده علوم ریاضی

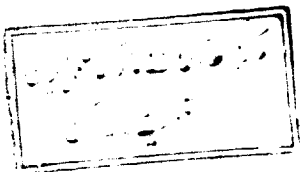
طبقه‌بندی داده‌ها با استفاده از مدل‌های ترتیبی

پایان نامه کارشناسی ارشد آمار
آیتین سعادت‌ملی

۱۹۹۹

استاد راهنما
دکتر ایوب ساعی

۱۳۸۰ / ۱ / ۱۰



۱۳۷۹

۳۲۷۹۹



دانشگاه صنعتی اصفهان
دانشکده علوم ریاضی

پایان نامه کارشناسی ارشد رشته آمار خانم آیتین سعادت ملی
تحت عنوان

طبقه‌بندی داده‌ها با استفاده از مدل‌های ترتیبی

در تاریخ ۷۹/۷/۲۵ توسط کمیته تخصصی زیر مورد بررسی و تصویب نهائی قرار گرفت.

دکتر ایوب ساعی

۱- استاد راهنمای پایان نامه

دکتر علی زینل همدانی

۲- استاد مشاور پایان نامه

دکتر مجتبی گنجعلی

۳- استاد داور ۱

دکتر محمد صالحی

۴- استاد داور ۲

دکتر امیر نادری

سرپرست تحصیلات تکمیلی دانشگاه

هذا من فضل ربي

حمد و سپاس خدای بزرگ و مهربان که همواره مرا مورد لطف و عنایت خویش قرار داده و توفیق انجام این پایان نامه را به من عطا فرمود.

بر خود لازم می دانم از کلیه عزیزان و سرورانی که مرا در انجام این پایان نامه یاری فرمودند کمال تشکر و قدردانی را بنمایم.

تشکر و سپاس ویژه از

جناب دکتر ایوب ساعی استاد راهنمای ارجمندم که با راهنماییهای دلسوزانه و نیز جدیت و پشتکار بسیار، بنده را در ارائه این پایان نامه یاری فرمودند.

جناب دکتر علی همدانی استاد مشاور که کمال همکاری و مساعدت را فرمودند.

جناب دکتر گنجعلی استاد داور ارجمند که زحمت بازخوانی و داوری این پایان نامه را متقبل شدند.

جناب دکتر محمد صالحی که به عنوان استاد داور دانشکده در جلسه دفاع حضور یافته و تذکرات جالب و مفیدی را مطرح کردند و با پیشنهادات ظریف و بجا زمینه هر چه بهتر شدن پایان نامه را فراهم فرمودند.

جا دارد از کلیه اساتید دانشکده ریاضی دانشگاه صنعتی به ویژه دکتر پارسسیان استاد گرانقدر و همچنین از کلیه اساتید گروه آمار دانشگاه اصفهان به ویژه دکتر علامتساز و به خصوص از خانواده خوب و فداکارم کمال تشکر و قدردانی را بنمایم.

کلیه حقوق مادی مترتب بر نتایج مطالعات،
ابتکارات و نوآوریهای ناشی از تحقیق موضوع
این پایان نامه (رساله) متعلق به دانشگاه صنعتی
اصفهان است.

تقدیر نام

پدر و مادر فداکاره

خواهر و برادر عزیزه

فهرست مطالب

صفحه	عنوان
شش	فهرست مطالب
۱	چکیده
۲	فصل صفر - مقدمه
۴	فصل یک - معرفی مدل‌هایی برای طبقه بندی داده‌ها
۵	۲-۱ - مدل‌های عددی
۹	۳-۱ - مدل رگرسیون لجستیک چندگانه
۱۴	۴-۱ - مدل‌های تشخیصی
۲۱	۵-۱ - مدل‌های بخت تجمعی
۲۶	۶-۱ - مدل نسبت دنباله‌ای
۲۹	۷-۱ - فرایند طبقه بندی
۳۵	۸-۱ - تحلیل ممیزی بین گروه‌های مرتب شده
۳۵	۹-۱ - تفاوت بین رگرسیون لجستیک و تحلیل ممیزی
۴۱	۱۰-۱ - نرخ خطاها
۴۸	فصل دوم - روش‌های شبیه سازی
۴۹	۱-۲ - متغیرهای کمکی دوتائی مستقل
۵۱	۲-۲ - متغیرهای کمکی نرمال چند متغیره

۵۱ ۳-۲ - داده‌های PO برای مقایسه فرایندهای PO و ML و AP

۵۴ ۴-۲ - داده‌های CR برای مقایسه فرایندهای ML و OL و ND

۵۷ فصل سوم - تحلیل ممیزی مدل‌های آستانه‌ای با اثرات تصادفی

۵۷ ۱-۳ - مقدمه

۵۸ ۲-۳ - مدل‌ها و نمادها

۷۰ ۳-۳ - برآورد

۷۳ ۴-۳ - کاربرد در بررسی بیماری‌های پوستی

۷۷ ۵-۳ - روش‌های شبیه‌سازی

۷۸ فصل چهارم - جدول‌ها و نتایج

چکیده

کمبل در سال ۱۹۹۱ یک روش شبیه سازی برای برازش مدل به داده‌های ترتیبی (Ordinal) انجام داد. نتایج شبیه سازی کمبل نشان داد که وقتی هدف اصلی طبقه بندی کردن داده‌ها باشد مدل‌های ترتیبی هیچ مزیتی بر مدل‌های غیر ترتیبی ندارد. کمبل در مقاله خود شبیه سازی را بدون در نظر گرفتن انواع مختلف داده‌ها و خواص آنها انجام داده است. در این تحقیق ما موضوع را با دیدگاه دیگری بررسی می‌کنیم. در این نوشتار با توجه به نوع گسسته یا پیوسته بودن متغیرهای کمکی مدل‌های مختلفی را مورد بررسی قرار داده و انواع خطاهائی را که از این طبقه بندی حاصل می‌شود به دست می‌آوریم. همچنین نشان داده خواهد شد که مدل‌های ترتیبی نسبت به مدل‌های غیر ترتیبی بهتر عمل کنند. در آخر طبقه بندی کردن داده‌ها را با مدل‌هایی با اثرات آمیخته مورد بررسی قرار می‌دهیم.

فصل صفر

مقدمه:

در دهه 60 میلادی ابتدا واکرو دانکن^(۱) مدل‌های مرتب شده را مورد توجه قرار دادند آنها مدل‌های بخت متناسب را مورد توجه قرار دادند پس مکالگ^(۲) (۱۹۸۰) به بسط این مدل‌ها که توسط واکرو دانکن مطرح شده بود پرداخت و آن را به فرم مدل‌های لجستیک و مقدارفرین ماکزیمال و مقدار فرین مینمال بسط داد و به بررسی آنها پرداخت.

در دهه ۷۰ ایشنیز و آوری^(۳) مدل‌های تحلیلی ممیزی را مورد بررسی قرار دادند که در سنوات بعد از آن به مدل‌های تحلیلی ممیزی نرمال با چندین گروه و سپس فرض با عدم همگنی، ماتریس واریانس کواریانس آن را بسط داده‌اند که این مطلب اخیر از حوصله این مقاله خارج است.

در سال ۱۹۸۱ فیلیپ و آندرسن^(۴) به بررسی رگرسیون مدل‌های تحلیلی ممیزی پرداختند و قواعدی برای طبقه بندی آنها ارائه کردند.

در سال ۱۹۸۴ میلادی آندرسن^(۵) به بررسی رگرسیون متغیرهای طبقه بندی شده مرتب شده پرداخت.

1- Walker fDuncan

2- Mccullagh

3- Eisenbers f Arery

4- Philips f Anderson

5- Anderson

بعد از آن ربشی و پولک^(۱) به بررسی رگرسیون اینگونه متغیرها در صورت دوتائی بودند همت گماردند و آرمسترانگ و اسلوان^(۲) مدل‌های نسبت متناسب را در سال ۱۹۸۹ ابداع کردند و در آن مسائلی را که مربوط به اپیدمیولوژی بود بررسی کردند. در نهایت کمبل^(۳) به بررسی خطاهای مربوط به مدل‌های مختلف پرداخت. البته او تنها خطای نوع اول که ترتیب را در نظر نمی‌گیرد مورد توجه قرار داد. در مقاله حاضر سه خطای مختلف را که خطای نوع اول آن همان است که کمبل در نظر گرفته بود و ۲ خطای نوع دوم و سوم را مورد بررسی قرار داده است. خطای نوع دوم و سوم خطاهای مربوط به متغیرها را در حالت ترتیبی یعنی زمانی که ترتیب مهم است در نظر گرفته‌ایم. در فصل دوم به روشهای شبیه سازی در حالتی که متغیرهای کمکی دوتائی یا نرمال چند متغیره است پرداخته‌ایم. در فصل سوم به بررسی همین مدلها در حالتی که متغیرهای تصادفی نیز در مدل می‌باشد پرداخته‌ایم.

در فصل چهارم نتایج کامپیوتری و تحلیل آنها آمده است. تمام محاسبات با استفاده از نرم‌افزار SAS و زیر برنامه PROC IML در آن انجام شده است.

1- Ashby f Pocok

2- Armstrong f sloan

3- cambell

فصل اول

معرفی مدل‌هایی برای طبقه بندی کردن

داده‌ها

(۱-۱) مقدمه

متغیرهای پاسخ گسسته که بیش از دو طبقه داشته باشند اغلب ترتیبی می‌باشند. بدین معنی که پیشامدهای توصیف شده توسط اعداد طبقه ۱، ۲، ...، J را می‌توان همچون ترتیب در نظر گرفت. فرض کنید X یک بردار P بعدی از متغیرهای پیش‌بینی کننده باشد و Y یک عدد مطابق با گروه عضویت باشد. فرض می‌کنیم که y ، J مقدار ۱، ۲، ...، J را می‌گیرد و به علاوه گروه‌های $Y = 1, \dots, J$ را به ترتیب بهترین (نرمال) و بدترین (بیمارترین و زیان‌آورترین) خروجی‌ها می‌نامیم. ترتیب پذیری Y به توزیع متغیرهای کمکی X منعکس می‌شود با میانگین آنها باید هم خطی چندگانه برای ترتیب پذیری مشاهدات داشته باشد.

ارتباط X, Y ممکن است با مدل‌های رگرسیون لجستیک بررسی شود. وقتی که $J=2$ است Y یک متغیر دوگروهی است؛ گروه‌ها خوب / بد یا قبول / رد می‌باشد و مدل رگرسیون لجستیک دوتائی (MLR) قابل کاربردی است. در عمل $J>2$ است یک مدل رگرسیون لجستیک چندگروهی (MLR) خواهیم

داشت. یک مثال سیستمهای چند گروهی بد / متوسط / خوب می باشد.

انگل^(۱) در سال (۱۹۸۸) گروهبندیهای زیر را بکار برد.

۱- مدل‌های عددی: این مدلها شامل تحلیل ممیزی نرمال (ND) و لجستیک چندگانه (ML) که در فصول بعدی شرح بیشتری درباره آن می دهیم می باشد.

۲- مدل‌های تشخیصی: این مدلها حالت‌هایی که در آنها روی گروه‌هایی که کارشناس تحلیل کرده است، شک و تردید وجود دارد را نشان می دهد. مثال این مدلها مدل لجستیک ترتیبی است.

۳- مدل‌های بخت تجمعی: این مدلها فرض می کند که پاسخهای ترتیبی به یک متغیر پیوسته و غیر قابل مشاهده وابسته است و این متغیر پنهان دارای یک توزیع خاص در یکی از فواصل متصل بهم که بانقاط بریدگی معلوم از هم جدا شده است قرار می گیرند و با حرکت از یک فاصله به فاصله دیگر متغیر پاسخ تغییر می کند.

یک مثال متغیر پنهان مطابق با (سرعت) ملایم / متوسط / سریع می باشد. مدل‌های بخت متناسب (PO) و نسبت دنباله‌ای (CR) دو زیر گروه این نوع مدلها می باشند.

در این ارتباط داریم:

$$\pi_{ij} = P(Y = j | X_i)$$

بطوریکه در آن π_{ij} از طریق قضیه بیز با عبارت‌هایی برای

$$\pi_j = P(Y = j), \quad P(X_i | Y = j) \text{ بدست می آید.}$$

گروه مرجع را $\lambda = 1$ زدر نظر می گیریم هدف از همه فرایندهایی که در زیر توضیح می دهیم برآورد تابع چند

جمله‌ای (۱-۱-۱) $\lambda_j(x_i) = \alpha_j + \beta_j' \times i$ می باشد.

(۱-۲) مدل‌های عددی:

۱- مدل‌های تحلیل ممیزی (ND):

این مدل توسط آوری^(۲) و ایشینبیز^(۳) (۱۹۷۵) پیشنهاد شد.

هدف از این تحلیل ممیزی چند متغیره این است که بر اساس اطلاعات قدرت جداسازی نمونه‌ها و یا صفات را بدست آورده و سپس برای مشاهدات جدید توانائی بهترین تخصیص را داشته باشیم. روش تحلیل ممیزی را یک روش جستجوگر وکشف کننده می‌نامیم.

اهداف استفاده از روش تحلیل ممیزی و طبقه بندی عبارتند از :

۱- تعیین تشخیص دهنده هائی که مقادیر عددی آنها برای مشاهدات چنان باشد که تا حد امکان جدا کردن آنها را ممکن سازد.

۲- مرتب نمودن مشاهدات و یا موضوعهای مورد مطالعه به چند گروه از قبل تعیین شده به طوریکه یک مشاهده را بتوان به صورت بهینه (یعنی حداقل خطا) به گروهی اختصاص داد.

فرمول بندی مسئله تحلیل ممیزی به این صورت است: فرض کنید J جمعیت متمایز π_1, \dots, π_J داریم. می‌خواهیم فردی را با مشاهده $x = (x_1, \dots, x_p)'$ با P خصوصیت مختلف که فرد را مشخص می‌سازد به یکی از جمعیت‌های π_1, \dots, π_J رده بندی کنیم.

وقتی گروه افراد را در نظر می‌گیریم اساساً فرض می‌کنیم که گروه به عنوان کل فقط به یکی از J جمعیت مفروض متعلق است. علاوه بر این فرض خواهیم کرد که هر یک از π_j ها را می‌توان بوسیله میانگینهای تابع توزیع F که به صورت یک بردار متغیر تصادفی $X = (X_1, \dots, X_p)$ می‌باشد تعیین کرد به قسمی که مؤلفه‌های این بردار تصادفی اندازه‌های تصادفی برای P خصوصیت مختلف را نشان می‌دهند.

تمام فضای P بعدی مقادیر X را به E^P نشان می‌دهیم. اکنون می‌خواهیم برای تقسیم E^P به J ناحیه مجزای R_1, \dots, R_J قاعده‌ای تنظیم کنیم به گونه‌ای که اگر X در R_j قرار گیرد فرد را به گروه π_j اختصاص می‌دهیم، بدیهی است که در بکار بردن چنین قاعده‌ای ممکن است با اشتباه رده بندی کردن یک فرد به π_i برای $i \neq j$ وقتی واقعاً فرد به π_j متعلق است مرتکب خطا شویم.

چون حالت K جمعیت نرمال چند متغیره حالت کلی تر از ۲ جمعیت نرمال چند متغیره است

ابتدا همان حالت $K = ۲$ را بررسی می‌کنیم.

۷

فرض کنید که $z = y$ داده شده باشد، X دارای توزیع نرمال چند متغیره P بعدی با بردار میانگین μ_i و ماتریس واریانس کوواریانس Σ است. آنگاه تابع داده شد در (۱-۱-۱) را می‌خواهیم برآورد کنیم.

$$F_i(x) = (\sqrt{2\pi})^{-p/\lambda} |\Sigma|^{-1/\lambda} \exp \left[\frac{-1}{\lambda} (x - \mu_i)' \Sigma^{-1} (x - \mu_i) \right]$$

$$\frac{F_1(x)}{F_2(x)} = \exp \left[\frac{-1}{\lambda} (x - \mu_1)' \Sigma^{-1} (x - \mu_1) + \frac{1}{\lambda} (x - \mu_2)' \Sigma^{-1} (x - \mu_2) \right]$$

$$= \exp \left[\frac{-1}{\lambda} (\mu_1 - \mu_2)' \Sigma^{-1} (\mu_1 + \mu_2) + (\mu_1 - \mu_2)' \Sigma^{-1} x \right]$$

$$\ln \frac{F_1(x)}{F_2(x)} = (\mu_1 - \mu_2)' \Sigma^{-1} x - \frac{1}{\lambda} (\mu_1 - \mu_2)' \Sigma^{-1} (\mu_1 + \mu_2)$$

$$L(x) = (\mu_1 - \mu_2)' \Sigma^{-1} \left[x - \frac{1}{\lambda} (\mu_1 + \mu_2) \right]$$

$$L(x) = a_0 + a_1 x$$

اما از طرفی

پس با مساوی قرار دادن دو عبارت $L(x)$ با هم داریم:

$$a_0 = \frac{-1}{\lambda} (\mu_1 - \mu_2)' \Sigma^{-1} (\mu_1 + \mu_2)$$

$$a = a_1 = (\mu_1 - \mu_2)' \Sigma^{-1}$$

که عبارت اخیر یک بردار سطری می‌باشد.

از اینرو اگر

$$a_0 + a_1 x = (\mu_1 - \mu_2)' \Sigma^{-1} x - \frac{1}{\lambda} (\mu_1 - \mu_2)' \Sigma^{-1} (\mu_1 + \mu_2)$$

$$\Rightarrow \hat{a}_0 = \frac{-1}{\lambda} (\mu_1 - \mu_2)' \Sigma^{-1} (\mu_1 + \mu_2) \quad (1-2-1)$$

$$\hat{a}_1 = (\mu_1 - \mu_2)' \Sigma^{-1} \quad (1-2-2)$$