

دانشگاه پیام نور

دانشکده علوم پایه

پایان نامه کارشناسی ارشد آمار ریاضی

مدل‌بندی آماری با استفاده از الگوی نرمال چند متغیره آمیخته

توسط:

امیر انشکانی

استاد راهنما:

دکتر پرویز نصیری

استاد مشاور:

دکتر علی شادرخ

تابستان ۱۳۹۱

شماره:

تاریخ:

پیوست:



دانشگاه پیام نور
دانشگاه پیام نور استان تهران
العلم بكل لؤلؤك الفرج والعاقبة والفضل



جمهوری اسلامی ایران
وزارت علوم، تحقیقات و فناوری

مرکز تهران شرق

صور تجلسه دفاع از پایان نامه کارشناسی ارشد

جلسه دفاع از پایان نامه کارشناسی ارشد آقای امیر انشکانی
دانشجوی رشته آمار ریاضی به شماره دانشجویی: ۸۷۰۰۰۰۲۶۰
تحت عنوان:

مدل بندی آماری با استفاده از الگوی نرمال چند متغیره آمیخته

جلسه دفاع با حضور داوران نامبرده ذیل در روز سه شنبه مورخ: ۹۱/۰۴/۲۷ ساعت: ۱۲-۱۱ در محل

تهران شرق برگزار شد. و پس از بررسی پایان نامه مذکور با نمره به عدد ۱۸.۹

به حروف هجده و نهم درصد و با درجه ارزشیابی کاملمورد قبول واقع شد نشد

ردیف	نام و نام خانوادگی	هیات داوران	مرتبه دانشگاهی	دانشگاه/موسسه	امضاء
۱	دکتر پرویز نصیری	استاد راهنما	دانشیار	پیام نور	
۲	دکتر علی شادرخ	استاد مشاور	استادیار	پیام نور	
۳	دکتر احسان جمالی	استاد داور	استادیار	سازمان سنجش آموزش کشور	
۴	دکتر مسعود یار محمدی	نماینده علمی گروه/ نماینده تحصیلات تکمیلی	دانشیار	پیام نور	

تهران ، خیابان کریمخان
زند ، خیابان استاد نجات
الهی ، خیابان شهید فلاح
پور ، پلاک ۲۷ مرکز
تهران شرق

تلفن: ۸۸۹۱۳۴۷۵

دورنگار: ۸۸۹۴۸۹۸۴

Tshargh.Tpnu.ac.ir

Tshargh@Tpnu.ac.ir

اگر شایسته تقدیم باشد

به روح بلند پرواز پدرم

که مشوق من بوده، هست و خواهد بود.

قدردانی

به نام یگانه پروردگار بزرگ

شکرا ایزد منان را که به من توانایی تفکر و عمل کردن آموخت. برخورد لازم می‌دانم که پس از پایان دوره کارشناسی ارشد از تمامی کسانی که در طول این دوره در کنار من بوده‌اند، به خصوص اعضای خانواده‌ام، تقدیر و تشکر نموده و درود و صلوات نثار روح پدر مرحومم کنم. همچنین از آقای دکتر پرویز نصیری که زحمت راهنمایی این پایان نامه و جناب آقای دکتر علی شادرخ که مشاوره بنده را قبول نمودند، به خاطر تمامی زحماتشان تشکر می‌کنم. امید که این تحقیق سهمی اندک در اعتلای علمی کشور داشته باشد.

امیرانشکانی

تهران

تابستان ۱۳۹۱

چکیده

در بسیاری از مسائل کاربردی تشخیص مدل مناسب برای توزیع صفت جامعه مورد بررسی از اهمیت ویژه‌ای برخوردار است. کاربرد و استفاده از توزیع‌های آمیخته به عنوان یکی از روش‌های داده‌کاوی سابقه بسیار طولانی در اغلب زمینه‌های علمی از جمله پزشکی، کشاورزی، هواشناسی، بازاریابی، مدیریت و غیره دارد. مزیت استفاده از اینگونه توزیع‌ها زمانی که داده‌ها دارای پراکندگی زیاد و حتی گمشده باشند، بدلیل انعطاف پذیری فوق‌العاده آنها به خوبی معلوم است. این نوع مسائل آماری اولین بار مورد توجه فلر قرار گرفت و سپس در مقالات متعددی روی این موضوع بحث و بررسی شد. توزیع‌های آمیخته را می‌توان برپایه بسیاری از توزیع‌های شناخته شده مثل یکنواخت، نمایی، نرمال، پواسن، دوجمله‌ای و غیره مدل بندی کرد. هدف از هر مدل‌بندی آماری برآورد پارامترهاست، برای برآورد پارامترها روشهای گوناگونی از جمله روش گشتاورها، حداقل مربعات خطا، روش ماکسیمم درستنمایی، روش مونت کارلو، الگوریتم EM و الگوریتم $MCMC$ وجود دارد.

واژه‌های کلیدی : نرمال چند متغیره آمیخته، شناسایی پذیری، هم واریانس، ناهم واریانس،

تابع درستنمایی

KeyWords: Mixture of multivariate normal, Identifiability, Homoscedastic, Heteroscedasticity, Likelihood function

فهرست مندرجات

۱	توزیع نرمال و تعمیم چند متغیره آن	۱
۱	۱.۱ مقدمه	۱
۱	۲.۱ توزیع نرمال یک متغیره	۱
۴	۱.۲.۱ توزیع های نمونه گیری	۴
۷	۲.۲.۱ گشتاورهای توزیع نرمال	۷
۹	۳.۱ توزیع نرمال چند متغیره	۹
۱۳	۱.۳.۱ گشتاورهای توزیع نرمال چند متغیره	۱۳
۱۵	۲ مدل های آمیخته	۱۵

فهرست مندرجات

۱۵	مقدمه	۱.۲
۱۶	تاریخچه مدل‌های آمیخته	۲.۲
۲۱	شناسایی‌پذیری مدل‌های آمیخته	۳.۲
۲۳	برآورد پارامترها	۴.۲
۳۶	انتخاب مقادیر اولیه الگوریتم EM در توزیع‌های آمیخته	۱.۴.۲

۳ مدل‌های نرمال و نرمال چند متغیره آمیخته

۳۸	مقدمه	۱.۳
۴۱	ترکیبات هم واریانس و ناهم واریانس	۲.۳
۴۶	ترکیبات هم واریانس	۱.۲.۳
۴۸	ترکیبات ناهم واریانس	۲.۲.۳

۴ آزمون‌های آمیخته

۵۱

۵۱ مقدمه ۱.۴

۵۳ آزمون‌های نسبت درست‌نمایی آمیخته ۲.۴

۵ مثال کاربردی و مطالعات بیشتر ۶۳

۶۳ مقدمه ۱.۵

۶۴ مدل‌بندی شاخص‌های هواشناسی شهر تهران ۲.۵

۶۵ حساسیت PTC ۳.۵

۶۸ ادامه مثال ۱.۴.۲ ۴.۵

۷۰ ارزیابی و مقایسه‌ی توابع تاوان (۸.۲.۴) و (۹.۲.۴) ۵.۵

۷۶ نتیجه‌گیری ۶.۵

فهرست مندرجات

۷۸ A برخی از اصطلاحات استفاده شده در متن پایان نامه

۸۲ B واژه‌نامه انگلیسی به فارسی

۸۶ C کتابنامه

۹۴ D ABSTRACT

لیست اشکال

۱.۲.۱ تابع چگالی توزیع نرمال $X \sim N(\mu, \sigma^2)$ ۲

۱.۲.۲ آمیخته گوسی (سمت چپ: ۲ تابع چگالی گوسی با

$\mu_1 = 0, \sigma_1^2 = 1$ و $\mu_2 = 3, \sigma_2^2 = 0.5$; سمت راست: آمیخته همان دو

توزیع با $(\pi_1 = 0.4, \pi_2 = 0.6)$ ۱۸

۱.۲.۳ نمودار چگالی نرمال ۲ متغیره با میانگین‌های (۱ و ۱) و (۱- و ۱-) .. ۴۴

۲.۲.۳ نمودار چگالی نرمال آمیخته ۲ متغیره با واریانس‌های برابر با ۱،

نسبت‌های برابر و $\mu_1 = 0, \mu_2 = \Delta$: (b) $\Delta = 3$, (c) $\Delta = 4$, (d)

۴۷ $\Delta = 2$, (a) $\Delta = 1$

- ۳.۲.۳ نمودار چگالی نرمال آمیخته ۲ متغیره با واریانس‌های برابر
با ۱، نسبت‌های ۰.۲۵ و ۰.۷۵ و میانگین‌های $\mu_1 = 0$, $\mu_2 = \Delta$
۴۷ (d) : $\Delta = 4$, (c) : $\Delta = 3$, (b) : $\Delta = 2$, (a) : $\Delta = 1$
- ۴.۲.۳ نمودار چگالی‌های نرمال آمیخته (منبع: *Marron and Wand*) ۴۹
- ۵.۲.۳ ادامه شکل ۴.۲.۳ ۵۰
- ۱.۲.۴ مقایسه نموداری توابع تاوان (۸.۲.۴) و (۹.۲.۴) ۶۱
- ۱.۵.۵ مقایسه توابع (۸.۴.۲) و (۹.۴.۲) با استفاده از نمودار خطی ۷۳

لیست جداول

- ۱.۳.۵ برآزش مدل آمیخته برای ۳ مجموعه اعداد (اعداد داخل پرانتز انحراف معیار هستند) ۶۷
- ۲.۴.۵ ۶۹
- ۳.۵.۵ توزیع‌های آمیخته پواسن و نرمال با میانگین‌های برابر تحت فرض‌های H_0 و H_1 ۷۱
- ۴.۵.۵ مقادیر برآورد شده و میانگین توان دوم خطای برآوردگرها به همراه آماره آزمون و توان آزمون نسبت درست‌نمایی اصلاح شده، به ازای توابع توان (۸.۲.۴) و (۹.۲.۴) ۷۲

لیست جداول

۵.۵.۵	میانگین توان دوم خطای برآوردگرهای بیزی و بسامدگرایانه در توزیع
۷۴	پواسن
۶.۵.۵	میانگین توان دوم خطای برآوردگرهای بیزی و بسامدگرایانه در توزیع
۷۵	نرمال
۷.۵.۵	مقادیر معیارهای نیکویی برازش آکائیک و گیبز
۷۶	

مقدمه

الگوی نرمال بدلیل برخورداری از برخی ویژگی‌های ذاتی و سهولتی که معمولاً در استنباط فراهم می‌آورد، از جایگاه ویژه‌ای در استنباط آماری برخوردار است. در سال‌های اخیر تلاش‌های گسترده‌ای برای توسعه الگوهای جامع‌تر از نرمال که علاوه بر ویژگی‌های مطلوب الگوی نرمال قابلیت بیشتری در مدل‌بندی داده‌های متنوع داشته باشند، صورت پذیرفته است. اما با وجود موفقیت آمیز بودن بخشی از این تلاش‌ها الگوی نرمال هنوز هم بدلیل پشتوانه استنباطی قابل توجهی که در ادبیات آماری دارد، دارای جایگاه ممتازی است. از طرف دیگر موارد بسیار زیادی وجود دارد که مشاهدات متعلق به چند الگوی یکسان اما با پارامترهای متفاوت هستند، چنین الگوهای را آمیخته گویند.

تحقیق و بررسی بر روی مدل‌های خاص آمیخته چند سالی است که در بین آماردانان رواج پیدا کرده است، از این رو پیدا کردن رابطه اینگونه توزیع‌ها با سایر موضوعات علمی از اهمیت ویژه‌ای برخوردار است. این در حالی است که اولین بار کارل پیرسن در سال ۱۸۹۴ اینگونه توزیع‌ها را معرفی نمود. بدنبال وی فنگ و مک کلاچ در سال ۱۹۹۲ شناسایی پذیری مدل‌های آمیخته را مطرح کردند. مک لاکلان و کریشنان (۱۹۹۴) مطالعات خود را بر الگوریتم EM

متمرکز نمودند و مگنوس و بولدا (۲۰۰۸) با ادغام الگوریتم EM و روش بوت استرپ، روشی را ابداع کردند که در نهایت منجر به افزایش سرعت همگرایی در داده‌ها و کاهش زمان و هزینه انجام محاسبات شد.

هدف اصلی در این پایان نامه بررسی مدل‌های چند متغیره نرمال آمیخته می‌باشد. در این تحقیق و در فصل اول مدل‌های نرمال یک متغیره و چند متغیره مورد مطالعه قرار گرفته‌اند. فصل دوم مربوط به مدل‌های آمیخته و مسائل مربوط به آن است. در فصل سوم مدل‌های نرمال و نرمال چند متغیره آمیخته به همراه ویژگی‌های آن‌ها و در فصل چهارم آزمون‌های شناخته شده روی مدل‌های آمیخته بررسی شده است. در فصل پنجم و پایانی نیز مثالی عددی و شبیه سازی اینگونه توزیع‌ها ارائه شده است.

فصل ۱

توزیع نرمال و تعمیم چند متغیره آن

۱.۱ مقدمه

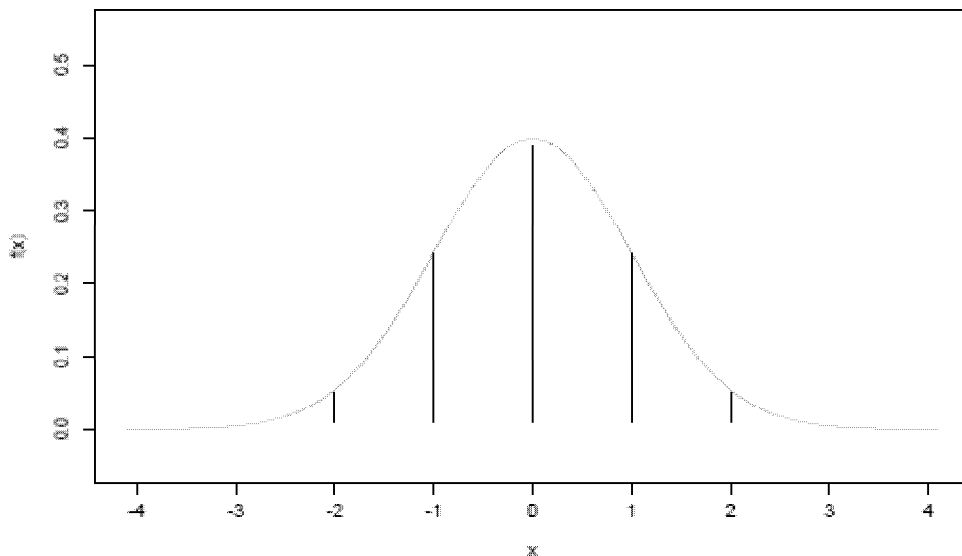
بدلیل اینکه در این پایان نامه از توزیع نرمال و ویژگی‌های آن زیاد استفاده می‌نماییم در این فصل به اجمال برخی از ویژگی‌های این توزیع مهم گردآوری شده است. در بخش دوم این فصل توزیع نرمال یک متغیره و در بخش سوم توزیع نرمال چند متغیره معرفی می‌گردد.

۲.۱ توزیع نرمال یک متغیره

تعریف ۱.۲.۱ متغیر تصادفی $-\infty < X < \infty$ دارای توزیع نرمال با میانگین $-\infty < \mu < \infty$ و واریانس $\sigma^2 > 0$ است $(X \sim N(\mu, \sigma^2))$ ، اگر تابع چگالی آن بصورت زیر تعریف شود.

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right], \quad (1.2.1)$$

همچنین نمودار تابع چگالی نرمال در شکل (۱.۵.۵) نشان داده شده است. برخی از



شکل ۱.۲.۱: تابع چگالی توزیع نرمال $X \sim N(\mu, \sigma^2)$

ویژگی‌های این توزیع در زیر فهرست شده است.

۱. برای یک فاصله غیر تهی (a, b) ، $f(x) \geq 0$ و

$$P(X \in (a, b)) = \text{area}_{(a,b)}.$$

که در آن $X \in (-\infty, +\infty)$.

۲. تابع چگالی f حول μ متقارن است، به عبارت دیگر $f(\mu + x) = f(\mu - x)$.

۳.

$$P(\mu - \sigma \leq X \leq \mu + \sigma) = 0/68$$

$$P(\mu - 2\sigma \leq X \leq \mu + 2\sigma) = 0/95$$

$$P(\mu - 3\sigma \leq X \leq \mu + 3\sigma) = 0/99.$$

۴. احتمال کوچکتر یا مساوی بودن متغیر تصادفی X با b یا سطح زیر منحنی نرمال از رابطه زیر بدست می آید.

$$F_X(b) = P(X \leq b) = \int_{-\infty}^b \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right] dx.$$

۵. اگر $X \sim N(\mu, \sigma^2)$

$$Z = \frac{X - \mu}{\sigma} \sim N(0, 1),$$

که Z را متغیر تصادفی نرمال استاندارد می نامند.

۶. اگر x_1, x_2, \dots, x_n نمونه‌ای تصادفی به حجم n از توزیع $N(\mu, \sigma^2)$ باشد. \bar{X} میانگین

نمونه دارای توزیع نرمال با میانگین μ و واریانس $\frac{\sigma^2}{n}$ است $(\bar{X} \sim N(\mu, \frac{\sigma^2}{n}))$.

نکته ۱ اگر x_1, x_2, \dots, x_n نمونه‌ای تصادفی به حجم n از توزیع $N(\mu, \sigma^2)$ باشد. به طوری که $-\infty < X < \infty$ ، $-\infty < \mu < \infty$ و $\sigma^2 > 0$ با استفاده از قضیه دسته بندی نیمن-فیشرف می توان اثبات نمود $(\sum_{i=1}^n x_i, \sum_{i=1}^n x_i^2)$ آماره بسنده مینیمال توام و کامل برای (μ, σ^2) است.

نکته ۲ برای یک نمونه تصادفی به حجم n از توزیع $N(\mu, \sigma^2)$ که $-\infty < X < \infty$ ، $-\infty < \mu < \infty$ و $\sigma^2 > 0$ ، \bar{X} و $\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$ برآوردگرهای ماکسیمم درست‌نمایی μ و σ^2 هستند.

نکته ۳ برای یک نمونه تصادفی به حجم n از توزیع $N(\mu, \sigma^2)$ که $-\infty < X < \infty$ ، $-\infty < \mu < \infty$ و $\sigma^2 > 0$ ، \bar{X} و $\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$ به ترتیب برآوردگرهای نارایب با کمترین واریانس برای μ و σ^2 هستند.

۱.۲.۱ توزیع‌های نمونه‌گیری

توزیع‌های نمونه‌گیری توزیع‌هایی هستند که براساس یک نمونه معمولاً n تایی از توزیع نرمال بدست می‌آیند. اینگونه توزیع‌ها با خواص جالب و ویژه خود در حالت مجانبی در مطالعات شبیه‌سازی کاربرد وسیعی دارند. توزیع‌های نمونه‌گیری بسیار زیادی وجود دارند اما مهمترین آنها توزیع‌های χ^2 ، T و F هستند. که در زیر به اجمال به برخی از ویژگی‌های آنها اشاره می‌شود.

تعریف ۲.۲.۱ در توزیع‌های نمونه‌گیری پارامتری به نام درجه آزادی وجود دارد، درجه آزادی تعداد متغیرهای تصادفی مستقل اینگونه توزیع‌ها است.

توزیع کای دو: ۱. اگر X داری توزیع نرمال استاندارد باشد، متغیر تصادفی X^2 دارای توزیع کای دو با یک درجه آزادی است. ($X \sim \chi^2_{(1)}$)

۲. در آمار از تابع چگالی زیر با عنوان توزیع کای دو با n درجه آزادی استفاده می‌شود.

$$f(x) = \frac{1}{2^{\frac{n}{2}} \Gamma(\frac{n}{2})} x^{\frac{n-2}{2}} e^{-\frac{x}{2}}; x > 0.$$

۳. امید ریاضی و واریانس متغیر تصادفی χ^2 با درجه آزادی ν به ترتیب ν و 2ν است.

۴. اگر $X_i \sim N(0, 1), i = 1, 2, \dots, n$ و مستقل از هم باشند، آنگاه $Y = \sum_{i=1}^n X_i^2$ دارای توزیع کای دو با n درجه آزادی است.

۵. اگر $X_i \sim \chi^2_{(i)}, i = 1, 2, \dots, n$ و مستقل از هم باشند، آنگاه $Y = \sum_{i=1}^n X_i$ دارای توزیع کای دو با $Y = \sum_{i=1}^n n_i$ درجه آزادی است.

۶. اگر \bar{X} و S^2 میانگین و واریانس نمونه تصادفی به حجم n و مستقل از هم باشند، متغیر تصادفی $\frac{(n-1)S^2}{\sigma^2} = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{\sigma^2}$ دارای توزیع کای دو با $(n-1)$ درجه آزادی است.

توزیع T : ۱. اگر Y, Z متغیرهای تصادفی مستقلی باشند، که Y دارای توزیع کای دو با ν

درجه آزادی و Z داری توزیع نرمال استاندارد باشند، آنگاه چگالی توزیع $T = \frac{Z}{\sqrt{\frac{Y}{\nu}}}$

دارای توزیع t با ν درجه آزادی است.