

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



دانشگاه اصفهان

دانشکده علوم

گروه آمار

پایان نامه ی کارشناسی ارشد رشته ی آمار گرایش آمار ریاضی

توزیع بتا - دو جمله ای: خواص و کاربردهای آن

استاد راهنما:

دکتر محمد بهرامی

پژوهشگر:

رامین جاهد حور

بهمن ماه ۱۳۸۹

کلیه حقوق مادی مترتب بر نتایج مطالعات،
ابتکارات و نوآوری های ناشی از تحقیق موضوع
این پایان نامه متعلق به دانشگاه اصفهان است.

سپاسگزاری

سپاس خدای را که مرا شوق و توان ادراک آموخت. با سپاس از کلیه اساتیدی که اندیشه امروز من، حاصل زحمات و تلاش‌های آنان است.

مراتب سپاس و تشکر خود را از استاد گرانقدر و ارجمندم جناب آقای دکتر محمد بهرامی که همواره با شکیبایی و درایت فراوان، مرا مستفیض راهنمایی‌های بی‌دریغ و ذی‌قیمت خود از لحاظ علمی و اخلاقی نمودند، ابراز می‌دارم. همچنین از زحمات بی‌دریغ اساتید داور جناب آقای دکتر هوشنگ طالبی و جناب آقای دکتر علی دولتی کمال تشکر را دارم.

از خانواده عزیزم مخصوصاً برادرم که تا این مرحله از زندگی همیشه حامی و پشتیبان من بوده و مادرم که صبر و ایثار را به من آموخت، نهایت تشکر و قدردانی را دارم. همچنین از دوستان و سایر عزیزان که موجد شوق فراگیری بیشتر در من می‌باشند، خالصانه و متواضعانه تشکر می‌نمایم.

رامین جاهد حور

بهمن ماه ۱۳۸۹

تقدیم به اولین معلمان و فرشتگان زندگیم

مادرم عزیز و فداکارم:

که نمی از زبودش، ایشار و گذشت کامل است

و نیمی دیگر عشق و محبت

برادران و خواهران مهربانم:

به پاس عاطفه سرشار و گرمای امید بخش وجودشان

و به پاس محبت های بی دریغ شان که هرگز فروکش نمی کنند.

و تقدیم به سرزمین مادریم

آذربایجان

چکیده

در نظریه آمار و احتمال، توزیع بتا-دوجمله‌ای خانواده‌ای از توزیع‌های احتمالی گسسته است که در بسیاری از تحقیقات کاربردی در زمینه‌هایی مانند علوم زیست‌شناسی، بهداشت، فیزیکی و اجتماعی مورد توجه محققان قرار گرفته است. انعطاف پذیری توزیع بتا-دوجمله‌ای باعث شده که این توزیع در برازش برخی مجموعه داده‌هایی که بیش پراکنش دارند سهم به سزایی ایفا کند، به طوری که در آمار بیزی، روش‌های بیز تجربی و آمار کلاسیک به عنوان توزیع دوجمله‌ای بیش پراکنده به صورت فراوان مورد استفاده قرار می‌گیرد. مطالعات فراوانی در مورد این توزیع انجام شده و روش‌های مختلفی برای برآورد کردن پارامترهای آن مورد بررسی قرار گرفته است. از آنجا که استفاده از روش‌های مناسب در برآورد پارامترهای این مدل دارای اهمیت ویژه‌ای می‌باشد، در این رساله، روش‌های مختلف برآوردیابی پارامترهای این مدل را مورد مطالعه قرار می‌دهیم. از جمله برآوردهای ماکسیمم درستنمایی را که به روش‌های عددی مانند روش نیوتن-رافسون به دست آمده است، از نظر کارایی جانبی نسبی با سایر روش‌های دیگر مقایسه می‌کنیم. در بخش دیگری از رساله حاضر تحلیل بیزی مدل بتا-دوجمله‌ای را مورد بررسی قرار می‌دهیم و به مطالعه آزمون‌های مختلف آماری برای این توزیع می‌پردازیم. با توجه به نتایج آزمون‌ها، می‌بینیم که وقتی در آزمون‌های تفاوت حساسیت یا ترجیح بیش از یک منبع تغییرات وجود دارد، مدل دوجمله‌ای داده‌های دوجمله‌ای بیش پراکنده را نمی‌تواند به خوبی برازش کند، ولی مدل بتا-دوجمله‌ای یک مدل جایگزین مناسب برای این آزمون‌ها می‌باشد. در ادامه، تقریب توزیع بتا-دوجمله‌ای را توسط توزیع دوجمله‌ای و پواسن ارائه می‌دهیم.

در پایان به شرح توزیع بتا-دوجمله‌ای تعمیم یافته که به وسیله‌ی تابع فوق هندسی گاوسی تولید شده است می‌پردازیم که به صورت توزیع آمیخته‌ای از بتا تعمیم یافته و توزیع دوجمله‌ای بیان شده است. این توزیع آمیخته جدید برای برازش برخی مجموعه داده‌ها مورد استفاده قرار می‌گیرد. با ارائه دو مثال نشان می‌دهیم که این توزیع می‌تواند برازش داده‌ها بوسیله توزیع بتا-دوجمله‌ای به دست آمده است را بهبود بخشد.

واژه‌های کلیدی: بوت استرپ، بیش پراکنش، تقریب پواسن، تقریب دوجمله‌ای، توان، توزیع بتا-دوجمله‌ای، توزیع دوجمله‌ای، شبیه سازی، کارایی جانبی نسبی، ماکسیمم درستنمایی.

فهرست مطالب

صفحه	عنوان
	فصل اول : تعاریف و مفاهیم پایه
۱-۱-۱	مقدمه.....
۱-۲-۱	تاریخچه تحقیق.....
۱-۳-۱	روش‌های برآوردیابی.....
۱-۳-۱-۱	روش گشتاوری.....
۱-۳-۱-۲	روش ماکسیمم درست‌نمایی.....
۱-۴-۱	توزیع‌های کاربردی و توابع فوق هندسی گاوسی.....
۱-۴-۱-۱	توزیع بتا.....
۱-۴-۱-۲	توزیع دوجمله‌ای.....
۱-۴-۱-۳	توابع فوق هندسی گاوسی.....
۱-۵-۱	روش‌های مهم و تعاریف اولیه.....
۱-۵-۱-۱	روش بوت استرپ.....
۱-۵-۱-۲	روش برآورد جانشینی.....
۱-۵-۱-۳	شبهه سازی مونت کارلو.....
۱-۵-۱-۴	آزمون نیکویی برازش.....
۱-۵-۱-۵	تعاریف اولیه در آمار بیز.....
	فصل دوم : توزیع بتا-دوجمله‌ای و برآورد پارامترهای آن
۱-۲-۱	مقدمه.....
۱-۲-۲	معرفی توزیع بتا-دوجمله‌ای (BBD).....
۱-۲-۲-۱	میانگین، واریانس و معیارهای دیگر متغیرهای تصادفی BB.....
۱-۲-۲-۲	مدل بتا-دوجمله‌ای به عنوان مدل کیسه.....
۱-۲-۲-۳	تابع مولد گشتاور، مولد احتمال و تابع مشخصه BBD.....

۱۸	۲-۲-۴- شکل‌های BBD به ازاء مقادیر مختلف پارامترهای آن.....
۲۱	۲-۳- برآورد پارامترها.....
۲۱	۲-۳-۱- برآورد گشتاوری پارامترها.....
۲۲	۲-۳-۲- برآورد ماکسیمم درست‌نمایی پارامترها.....
۲۴	۲-۴- روش‌های دیگر برآورد پارامترها.....
۲۵	۲-۴-۱- برآوردگرها براساس میانگین و صفرها (میانگین-صفرها).....
۲۶	۲-۴-۲- برآوردگرها براساس دو گشتاور نمونه‌ای اول (دو گشتاوری).....
۲۶	۲-۴-۳- برآوردگرها براساس میانگین و نسبت یک‌ها به صفرها (یک گشتاور-یک احتمال).....
۲۷	۲-۴-۴- برآوردگرهای مینیمم خی-دو.....
۲۸	۲-۵- مقایسه کارایی‌های مجانبی نسبی.....
۳۲	۲-۶- تحلیل بیزی مدل بتا-دوجمله‌ای.....
۳۲	۲-۶-۱- رهیافت بیزی.....
۳۵	۲-۶-۲- روش بیزی برای پیشگویی تعداد وقوع رخدادها.....
۳۵	۲-۶-۳- مطالعات شبیه سازی.....
۴۵	۲-۷- نتیجه گیری.....

فصل سوم : آزمون‌های مختلف آماری برای توزیع بتا-دوجمله‌ای

۴۶	۳-۱- مقدمه.....
۴۷	۳-۲- آزمون‌هایی برای یک آزمایش.....
۴۸	۳-۳- آزمون‌هایی برای دو آزمایش مستقل.....
۴۸	۳-۴- آزمون‌هایی برای بیش از دو آزمایش مستقل.....
۴۹	۳-۵- آزمون نسبت درست‌نمایی.....
۴۹	۳-۶- آزمون نیکویی برازش.....
۵۰	۳-۷- مثال‌های عددی.....
۵۵	۳-۸- تعیین اندازه نمونه براساس دقت برآورد و توان آزمون.....

۳-۸-۱- معیار دقت	۵۵
۳-۸-۲- معیار توان	۵۷
۳-۹-۹- مطالعات شبیه سازی	۵۹
۳-۹-۱- شبیه سازی خطای نوع اول برای آزمون‌های دوجمله‌ای و BB	۶۰
۳-۹-۲- شبیه سازی توان آزمون BB	۶۰
۳-۹-۳- شبیه سازی توان‌های آزمون‌های BB	۶۱
۳-۱۰-۱- آزمون نیکویی برازش بوت استرپ مدل BB	۶۲
۳-۱۰-۱- رویکردهای آزمون نیکویی برازش	۶۳
۳-۱۱-۱- آزمون‌های توزیع تجربی مدل BB	۶۴
۳-۱۱-۱- آزمون GOF براساس EDF	۶۵
۳-۱۱-۲- مطالعات شبیه سازی	۶۸
۳-۱۲- نتیجه گیری	۷۶
فصل چهارم : تقریب توزیع بتا-دوجمله‌ای با توزیع‌های دوجمله‌ای و پواسن	
۴-۱- مقدمه	۷۷
۴-۲- تقریب دوجمله‌ای توزیع بتا-دوجمله‌ای	۷۸
۴-۲-۱- اتحاد اشتاین	۷۹
۴-۳- تقریب پواسن توزیع بتا-دوجمله‌ای	۸۲
۴-۳-۱- اتحاد اشتاین-چن	۸۲
فصل پنجم : توزیع بتا-دوجمله‌ای تعمیم یافته	
۵-۱- مقدمه	۸۵
۵-۲- توزیع تعمیم یافته بتا-دوجمله‌ای	۸۷
۵-۲-۱- پیدایش توزیع	۸۷
۵-۲-۲- توزیع بتا-دوجمله‌ای تعمیم یافته به عنوان دوجمله‌ای آمیخته	۸۹
۵-۳- خواص توزیع	۹۱

عنوان

صفحه

۹۱۵-۳-۱-افراز واریانس.....
۹۲۵-۳-۲- میانگین احتمال و همبستگی بین پیشامدها.....
۹۳۵-۴-۴-مثالها.....
۹۳۵-۴-۱- مثال مصرف روزانه الكل.....
۹۵۵-۴-۲- مثال نمرات قبولی دانشگاهی.....
۱۰۰ منابع و مأخذ.....

فهرست شکل‌ها

صفحه

عنوان

- شکل ۱-۲- شکل‌های مختلف توزیع BB به ازاء مقادیر مختلف پارامترهای آن..... ۱۹
- شکل ۲-۲- نمودارهای همگرایی برآورد ML پارامترهای α و β در ۱۰ تکرار فرآیند نیوتن-رافسون... ۲۴
- شکل ۱-۳- سطح دقت برآورد پارامتر برای مدل BB (برای $\mu = 0.7$ ، $\rho = 0.5$ و $\alpha = 0.1$)..... ۵۶
- شکل ۲-۳- سطح توان آزمون BB (برای $\alpha = 0.1$ ، $\mu_0 = 0.5$ ، $\mu_1 = 0.7$ و $\rho = 0.5$)..... ۵۸
- شکل ۳-۳- منحنی‌های توان یکسان (برای $\alpha = 0.1$ ، $\mu_0 = 0.5$ ، $\mu_1 = 0.7$ و $\rho = 0.5$)..... ۵۹
- شکل ۱-۴- تابع احتمال دو جمله ای و بتا-دو جمله ای براساس $\alpha = 0.4$ و $\beta = 10$ ۸۱
- شکل ۲-۴- تابع احتمال دو جمله ای و بتا دو جمله ای براساس $\alpha = 0.3$ و $\beta = 100$ ۸۱
- شکل ۳-۴- تابع احتمال بتا دو جمله ای و پواسن براساس $\alpha = 20$ و $\beta = 30$ ۸۴
- شکل ۳-۴- تابع احتمال بتا دو جمله ای و پواسن براساس $\alpha = 20$ و $\beta = 100$ ۸۴
- شکل ۱-۵- (الف) $PMF_{GBB_{10}}(0.8, 1.2, 1.2)$ (*) و $BB_{10}(0.8, 1.2)$ (ب) چگالی‌های $Gbeta_{10}(0.8, 1.2, 1.2)$ (-) و $beta(0.8, 1.2)$ ۸۹

فهرست جدول‌ها

صفحه	عنوان
۱۲	جدول ۱-۱- خانواده توزیع‌های مزدوج
۳۰	جدول ۱-۲- ARE برآوردگرها برای پارامترهای خانواده بتا - دو جمله‌ای وقتی $n=5$
۳۱	جدول ۲-۲- ARE برآوردگرها برای خانواده بتا - دو جمله‌ای وقتی $n=10$
۳۶	جدول ۳-۲- خلاصه آماره‌های اصلی برای سه نوع پیشین
۳۸	جدول ۴-۲- مقایسه برآورد با پیشین‌های مختلف با استفاده از داده‌های شبیه‌سازی شده با $n=20$ و $r=2$
۳۹	جدول ۵-۲- مقایسه برآورد با پیشین‌های مختلف با استفاده از داده‌های شبیه‌سازی شده با $n=20$ و $r=3$
۴۰	جدول ۶-۲- مقایسه برآورد با پیشین‌های مختلف با استفاده از داده‌های شبیه‌سازی شده با $n=20$ و $r=4$
۴۱	جدول ۷-۲- مقایسه برآورد با پیشین‌های مختلف با استفاده از داده‌های شبیه‌سازی شده با $n=30$ و $r=2$
۴۲	جدول ۸-۲- مقایسه برآورد با پیشین‌های مختلف با استفاده از داده‌های شبیه‌سازی شده با $n=30$ و $r=3$
۴۳	جدول ۹-۲- مقایسه برآورد با پیشین‌های مختلف با استفاده از داده‌های شبیه‌سازی شده با $n=30$ و $r=4$
۴۴	جدول ۱۰-۲- مقایسه برآوردگرهای فاصله‌ای متقارن و نامتقارن بیزی برای p با استفاده از داده شبیه‌سازیشده و پیشین‌های یکنواخت
۵۱	جدول ۱-۳- (الف) داده‌های مثال ۱ و ۳
۵۱	جدول ۱-۳- (ب) داده‌های مثال ۲

- جدول ۳-۱- (پ) داده‌های مثال ۴ ۵۱
- جدول ۳-۲- مقادیر اندازه نمونه‌های n و k (برای $\Delta = 0.1$ ، $\mu = 0.7$ و $\rho = 0.5$ و در سطح $\alpha = 0.1$ (در آزمون یکطرفه) ۵۷
- جدول ۳-۳- مقایسه توان‌های آزمون توزیع بتا-دوجمله‌ای (برای $n=12$ ، $k=18$ و $\mu=0.6$) ۶۱
- جدول ۳-۴- مقایسه توان‌های آزمون توزیع بتا-دوجمله‌ای (برای $n=12$ ، $k=18$ و $\mu=0.6$) ۶۲
- جدول ۳-۵- سطح تجربی آماره‌های منتخب براساس ۱۰۰۰ تکرار آزمایش مونت کارلو ۶۹
- جدول ۳-۶- توان‌های تجربی (بر حسب درصد) برای آماره‌های منتخب براساس ۱۰۰۰ تکرار مونت کارلو در $\alpha = 0.05$ ۷۰
- جدول ۳-۷- سطح تجربی آزمون پیشنهاد شده برای مدل بتا-دوجمله‌ای براساس ۵۰۰۰ تکرار مونت کارلو ۷۲
- جدول ۳-۸- توان تجربی آماره‌های منتخب براساس ۵۰۰۰ تکرار مونت کارلو در سطح $\alpha = 0.05$ ۷۴
- جدول ۵-۱- تعداد روزهای نوشیدن الکل ۸۶
- جدول ۵-۲- تعداد موارد هر دانشجوی، سال تحصیلی ۲۰۰۳-۲۰۰۴ ۸۷
- جدول ۵-۳- برآوردهای پارامترها و نیکویی برازش تعداد روزهای نوشیدن الکل (با خطاهای استاندارد در پرانتزها) ۹۴
- جدول ۵-۴- برآوردهای متوسط احتمال مصرف روزانه ($\hat{\mu}$)، ضریب همبستگی بین مصرف در روزهای مختلف ($\hat{\rho}$)، درصد واریانس تفاوت بین افراد ($\% (100\hat{\sigma}_{nr}^2 / \hat{\sigma}^2)$) و درصد واریانس شانس ($\% (100\hat{\sigma}_r^2 / \hat{\sigma}^2)$) ۹۵
- جدول ۵-۵- برآورد پارامترها و AIC برای تعداد درس‌های پاس شده، که بوسیله مدل‌های BB و GBB برازش شده‌اند ۹۷
- جدول ۵-۶- احتمال تجربی متوسط قبولی درس و برآورد شده ($\hat{\mu}_n, \mu_n$)، ضریب همبستگی برآورد شده بین قبولی در دروس مختلف ($\hat{\rho}_n$)، واریانس کل برآورد شده ($\hat{\sigma}_n^2$) و درصد متعلق به آن به علت تغییرپذیری بین افراد در نمونه ($\% (100\hat{\sigma}_{nr}^2 / \hat{\sigma}_n^2)$) و نسبت بین واریانس‌های تجربی و برآورد شده ($s_n^2 / \hat{\sigma}_n^2$) ۹۸

مخفف‌ها و نمادها

AIC.....	AkaikInformaition Criterion
ARE.....	Asymptotic Relative Efficiency
AW.....	Average Width of the Interval
BBD.....	Beta-Binomial Distribution
CV.....	Coefficient of Variation
d_{TV}	Total Variation Distance
EDF.....	Empirical-Distribution-Function
GBBD.....	Generalization of the Beta-Binomial Distribution
GHF.....	Gaussian Hypergeometric Functions
GOF.....	Goodness-of-Fit
INCL.....	Proportion of Intervals Including the True Value
K.....	Kurtosis
LR.....	Likelihood Ratio
MAD.....	Mean Absolute Deviations
MME.....	Method ofMomentEstimate
MLE.....	Maximum Likelihood Estimate
MSE.....	Mean Square Error
NBD.....	Negative Binomial Distribution
PMF.....	ProbabilityMass Function
RMSE.....	Root of Mean Square Errors
SD.....	Standard Deviation
SK.....	Skewnes

فصل اول

تعاریف و مفاهیم پایه

۱ - ۱ مقدمه

در این فصل به معرفی مفاهیم و تعاریفی می‌پردازیم که در فصل‌های بعد آنها را مورد استفاده قرار خواهیم داد. در بخش دوم، ابتدا به تاریخچه تحقیق و کارهایی که در این زمینه صورت گرفته می‌پردازیم. در بخش سوم روش‌های برآوردیابی گشتاوری و ماکسیمم درستنمایی را بیان می‌کنیم. سپس در بخش چهارم به معرفی برخی توزیع‌های مهم پرداخته و بعضی از ویژگی‌ها و مشخصه‌های آنها را شرح می‌دهیم. روش‌های بوت استرپ و برآورد جانشینی و شبیه‌سازی مونت کارلو رادر بخش پنجم معرفی کرده و در نهایت مفاهیم پایه‌ای در آمار بیز را شرح خواهیم داد.

۱-۲ تاریخچه تحقیق

در نظریه احتمال و آمار توزیع بتا - دوجمله‌ای خانواده‌ای از توزیع‌های احتمالی گسسته آمیخته است که در آمار بیزی، روش‌های بیز تجربی و آمار کلاسیک به عنوان توزیع دوجمله‌ای بیش پراکنده^۱ به طور گسترده‌ای مورد استفاده قرار می‌گیرد. این توزیع کاربردهای فراوانی در علوم اجتماعی، فیزیک و بهداشت دارد. مدل بتا-

¹Overdispersed Binomial Distribution

دوجمله‌ای اولین بار به صورت رسمی توسط اسکلام^۱ (۱۹۴۸) ارائه شد، با این وجود این ایده قبلاً توسط پیرسن (۱۹۲۵) در یک آزمایش تحقیقی تئوری نیز پیشنهاد شده بود.

تنوع و انعطاف‌پذیری مدل بتا - دو جمله‌ای باعث شد که از آن در آزمون هوش (هوین^۲، ۱۹۷۹؛ لرد^۳، ۱۹۶۵؛ ویلکاکس^۴، ۱۹۸۱)، آزمایش تاکسیکولوژی (ویلیامز^۵، ۱۹۷۵)، اپیدمیولوژی (گریفتس^۶، ۱۹۷۳)، نمایش رسانه‌ها (گرین^۷، ۱۹۷۰) و رفتار خریداران (مسی^۸ و همکاران، ۱۹۷۰) استفاده شود.

ویلیامز (۱۹۷۵) داده‌های مربوط به تأثیر دارویی خاص بر روی حیوانات آزمایشگاهی را مورد بررسی قرار داد، که در این آزمایش متغیرهای پاسخ X_i تعداد توله‌هایی است که در دوره‌ی شیرخوارگی زنده مانده‌اند. بنابراین در این آزمایش فرض می‌کنیم که X_i متغیر تصادفی دو جمله‌ای با پارامتر P است. اسکورنیک^۹ (۱۹۹۰) در طراحی یک آزمایش روی بیماران غیر بارور، جمعیتی از زنان را که به صورت منظم تخمک‌گذاری نمی‌کردند، مورد بررسی قرار داد. برای هر زن تعداد دوره‌های تخمک‌گذاری X_i از شش دوره‌ی متوالی ثبت شد که استاندارد کلینیکی درمان می‌باشد. بنابراین هر X_i یک متغیر تصادفی دوجمله‌ای با پارامتر P است. در مسئله حیوانات آزمایشگاهی میزان بقا توله‌ها در یک زایمان، از یک زایمان به دیگری تغییر می‌کند. مشابهاً در مطالعه میزان تخمک‌گذاری در پاسخ به درمان می‌توان انتظار داشت که از یک زن به زن دیگر تغییر کند. بنابراین توزیع تعداد چرخه‌های موفق درون گروه تیمار در مقایسه با مدل دوجمله‌ای بیش پراکندگی دارد، که میزان تخمک‌گذاری در هر گروه تخمین زده شده است.

در چنین مواردی، توزیع بتا - دوجمله‌ای مدل انعطاف‌پذیری برای تغییر پذیری بین نمونه فراهم می‌کند. در این مورد، فرض می‌شود که احتمال تخمک‌گذاری هر زن در یک چرخه تنها به صورت بتا باشد و توزیع تعداد چرخه‌های موفق، به شرط آن احتمال دوجمله‌ای باشد.

در مطالعه اثر زایمان روی مدل بندی دز پاسخ در تراتولوژی، کوپر^{۱۰} و همکارانش (۱۹۸۶) نتیجه گرفتند که عدم بکارگیری تأثیر زایمان می‌تواند باعث کم برآوردی واریانس مرتبط با برآوردهای پارامتر باشد، به طوری که

¹Skellam

²Huynth

³Lord

⁴Wilcox

⁵Williams

⁶Griffiths

⁷Green

⁸Massy

⁹Skurnick

¹⁰Kupper

استفاده از درستنمایی دو جمله‌ای برای مدل‌بندی داده‌های تراتولوژی به نظر معقول نیست. آنها توزیع بتا-دوجمله‌ای را برای معرفی درجه تغییر همبستگی درون زایمان پیشنهاد کردند.

پُل^۱ (۱۹۸۲) در تحلیل نسبت‌های جنین تأثیر پذیر در آزمایش‌های تراتولوژیکی مشاهده کرد که مدل بتا-دوجمله‌ای نسبت به مدل دوجمله‌ای بهتر عمل کرده است. پک^۲ (۱۹۸۶) نتیجه گرفت که بتا-دوجمله‌ای نسبت به مدل‌های معادل مانند دوجمله‌ای همبسته کوپر و هاسمن^۳ (۱۹۷۸) بهتر است. تارون^۴ (۱۹۸۲) برای بسیاری از انواع انواع غده‌ها مشاهده کرد که نرخ کنترل عمر غده‌ها بسیار پذیرتر از آن هستند که فرض کنیم دارای توزیع دوجمله‌ای باشند، و توزیع بتا-دوجمله‌ای را برای نرخ عمر غده‌ها برآزش داد. شوکرز^۵ (۲۰۰۳) از توزیع بتا-دوجمله‌ای در ارزیابی کارایی دستگاه تشخیص هویت بیومتریکی استفاده کرد. لین و چو^۶ (۲۰۰۷) آزمون‌های تابع توزیع تجمعی را برای مدل بتا-دوجمله‌ای ارائه دادند.

تکنیک‌های مختلفی برای برآورد پارامترهای توزیع بتا-دوجمله‌ای ارائه شده است. اسکلام (۱۹۴۸) از برآوردهای ماکسیمم درستنمایی (MLE) که با کمک جداول تابع دایگاما^۷ به دست می‌آیند، استفاده کرد. گریفیس (۱۹۷۳) روشی برای به دست آوردن MLE توسعه داد که نیازی به استفاده از مقادیر تابع دایگاما ندارد، اما شامل فرایند تکراری برای حل معادلات غیر خطی است. نیسن-میر^۸ (۱۹۶۴) یک روش گرافیکی تکراری برای به دست آوردن MLE پیشنهاد کرد. ویلیامز (۱۹۸۲) روش ساده‌ای برای تغییر تحلیل لجستیک خطی استاندارد بیان کرد، که از برنامه کامپیوتری GLIM برای در نظر گرفتن تغییر پذیری دوجمله‌ای اضافی استفاده ویژه‌ای می‌کند. ویلیامز (۱۹۸۸) ارزیابی مرتبط با MLE را مطالعه کرد. بروکس^۹ (۱۹۸۴) روش تغییر یافته‌ای از GLIM پیشنهاد کرد که روش ساده‌ای است برای به دست آوردن تقریب آماری آزمون نسبت درستنمایی که توزیع بتا-دوجمله‌ای دارد. علاوه بر روش‌های MLE، چتفیلد و گودهارت^{۱۰} (۱۹۷۰) روش تکراری ساده برای محاسبه برآوردها بکار بردند که بر پایه معادله‌بندی میانگین نمونه و نسبت صفرهاست. شنتون^{۱۱} (۱۹۵۰) نشان داد

¹Paul

²Pack

³Kupper and Haseman

⁴Tarone

⁵Schuckers

⁶Lin and Chou

⁷Digamma Function

⁸Nissen-Meyer

⁹Brooks

¹⁰Chatfield and Goodhart

¹¹Shenton

که کارایی مجانبی نسبی روش گشتاورها معمولاً از ۷۰٪ بیشتر می‌باشد. همچنین انسکومب^۱ (۱۹۵۰) نشان داد که برای برازش توزیع دوجمله‌ای منفی (NB) که شکل معکوس J (فراوانی بزرگ صفرها) دارند روش میانگین و صفرها کارا تر است.

چون توزیع NB مورد منحصری از توزیع بتا - دوجمله‌ای است، چتفیلد و گودهارت (۱۹۷۰) روش میانگین و صفرها را برای برازش توزیع بتا-دوجمله‌ای که شکل معکوس J دارند، پیشنهاد کرد. ویکاکس (۱۹۷۹) با استفاده از تکنیک‌های مونت کارلو، تقریب نیوتن رافسون برآورد ماکسیمم درستنمایی توزیع بتا-دوجمله‌ای را با چندین روش ممکن دیگر مقایسه کرد.

۱ - ۳ روش‌های برآوردیابی

۱ - ۳ - ۱ روش گشتاوری

یکی از قدیمی‌ترین روش‌های برآوردیابی، روش برآورد گشتاوری است که در سال ۱۸۹۴ توسط آماردان مشهور کارل پیرسن معرفی شده است. روش برآورد گشتاوری که از آن با عنوان برآورد MM یاد می‌کنیم، عبارت است از دستورالعملی برای به دست آوردن برآوردگری به نام "برآورد گشتاوری" که از آن به اختصار MME یاد می‌کنیم. برآورد MM مبتنی بر بکارگیری گشتاورهای جمعیتی (توزیع) و گشتاورهای نمونه‌ای است.

برای روشن شدن موضوع، فرض کنید X_1, \dots, X_n یک نمونه‌ی تصادفی n تایی از توزیع $F_{\theta} \in \mathcal{F}$ باشد، به طوری که $\theta = (\theta_1, \dots, \theta_k) \in \Theta \subseteq R^k$. همچنین فرض کنید k گشتاور اول این توزیع، که به صورت توابعی از θ هستند، وجود داشته باشند. می‌دانیم که گشتاور I -م توزیع، در صورت وجود، به صورت زیر تعریف می‌شود:

$$\mu_r = \mu_r(\theta) = E_{\theta}(X_1^r) ; \quad r = 1, \dots, k .$$

اگر

$$M_r = \frac{1}{n} \sum_{i=1}^n X_i^r \quad , \quad r = 1, \dots, k$$

نمایانگر I -آمین گشتاور نمونه‌ای بر پایه نمونه تصادفی داده شده باشد، آنگاه برآورد MM پارامترهای مجهول $\theta_1, \dots, \theta_k$ براساس یک ایده ساده و از تشکیل و حل k معادله زیر حاصل خواهد شد

$$\mu_r = M_r \quad , \quad r = 1, \dots, k . \quad (1-1)$$

¹ Anscombe

معمولاً برآورد گشتاوری θ را با $\tilde{\theta}$ نمایش می‌دهند، این برآورد لزومی ندارد که منحصر به فرد باشد. بنابراین با توجه به آنچه که گفته شد، اگر $\theta = \varphi(\mu_1, \dots, \mu_r)$ باشد، آنگاه $\tilde{\theta} = \varphi(M_1, \dots, M_r)$ خواهد بود.

۱-۳-۲ روش ماکسیمم درست‌نمایی

این روش یکی از قدیمی‌ترین و پراهمیت‌ترین روش‌ها در نظریه برآوردهاست، که از آن به اختصار با عنوان برآورد ML یاد می‌کنیم. این روش اولین بار توسط گوس (۱۸۲۱) بکار گرفته شد و پس از آن به صورت گسترده‌تری در سال ۱۹۲۵ توسط فیشر در انتقاد از روش "برآورد گشتاوری" مورد استفاده قرار گرفت. روش ML عبارت است از دستورالعملی برای به دست آوردن برآوردگری به نام "برآوردگر ماکسیمم درست‌نمایی" که از آن به اختصار MLE یاد خواهیم کرد و مبتنی بر یک تابع آماری مهم به نام "تابع درست‌نمایی" است. فرض کنید $X = (X_1, \dots, X_n)$ بردار n متغیر تصادفی با تابع چگالی احتمال توأم $f_\theta(x)$ ، $\theta \in \Theta \subseteq R^k$ باشد.

تعریف ۱-۱ تابع درست‌نمایی برای هر مقدار داده شده $X = x$ ، تابع درست‌نمایی X را تابع چگالی احتمال توأم X ، یعنی $f_\theta(x)$ ، تعریف می‌کنیم که به صورت تابعی از θ در نظر گرفته می‌شود و آن را با نماد $L(\theta)$ نمایش می‌دهیم:

$$L(\theta) = f_\theta(x) = L(\theta; x).$$

ذکر چند نکته در ادامه ضروری به نظر می‌رسد.

نکته ۱-۱ تابع درست‌نمایی $L(\theta)$ لزوماً نسبت به θ مشتق پذیر نیست.

نکته ۱-۲ تابع درست‌نمایی $L(\theta)$ یک تابع چگالی احتمال نیست.

نکته ۱-۳ اگر X_1, \dots, X_n یک نمونه تصادفی n تایی از خانواده چگالی‌های $\{f_\theta(x): \theta \in \Theta\}$ باشد، آنگاه:

$$L(\theta) = \prod_{i=1}^n f_\theta(x_i)$$

تعریف ۱-۲ برآوردگر ماکسیمم درست‌نمایی اگر $\delta(x)$ برآوردگری برای θ باشد به طوری که:

$$i) \quad P_\theta(\delta(X) \in \Theta) = 1 \quad ; \quad \forall \theta \in \Theta$$

$$ii) \quad L(\delta(x)) \geq L(\theta) \quad ; \quad \forall \theta \in \Theta$$

آنگاه $\delta(x)$ به عنوان برآوردگر ماکسیمم درست‌نمایی θ تعریف می‌شود.

معمولاً بر آورد ما کسیم در ستمایی θ را با $\hat{\theta}$ نشان می دهند و لزومی ندارد که منحصر به فرد باشد. همچنین بر اساس تعریف $\hat{\theta}$ ، داریم:

$$L(\hat{\theta}) = \sup_{\theta \in \Theta} L(\theta)$$

به طور معادل اگر تعریف کنیم: $l(\theta) = \ln L(\theta)$ ، آنگاه

$$l(\hat{\theta}) = \sup_{\theta \in \Theta} l(\theta)$$

۱-۴ توزیع های کاربردی و توابع فوق هندسی گاوسی

در این بخش به مرور توزیع های بتا، دو جمله ای، و توابع فوق هندسی گاوسی^۱ می پردازیم و برخی از ویژگی های آنها را ارائه می کنیم.

۱-۴-۱ توزیع بتا

اگر متغیر تصادفی X دارای تابع چگالی زیر باشد

$$f_{\alpha, \beta}(x) = \frac{1}{B(\alpha, \beta)} x^{\alpha-1} (1-x)^{\beta-1}; \quad 0 < x < 1, \quad \alpha, \beta > 0 \quad (2-1)$$

آنگاه X دارای توزیع بتا با پارامترهای α و β است و آن را با نماد $X \sim \text{Beta}(\alpha, \beta)$ نشان می دهیم، که در

$$B(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha + \beta)}, \quad \text{داریم (2-1)}$$

در این توزیع $E(X) = \frac{\alpha}{\alpha + \beta}$ و $\text{Var}(X) = \frac{\alpha\beta}{(\alpha + \beta + 1)(\alpha + \beta)^2}$ است.

۱-۴-۲ توزیع دو جمله ای

هرگاه X تعداد موفقیت ها در n بار تکرار مستقل یک آزمایش بر نولی با احتمال موفقیت P در هر آزمایش

باشد ($0 < p < 1, n > 0$) آنگاه X را یک متغیر تصادفی دو جمله ای گوئیم و آن را با نماد $\text{Bin}(n, p)$

نشان می دهیم. تابع چگالی احتمال آن به این صورت است

$$f_X(x) = \binom{n}{x} p^x (1-p)^{n-x}, \quad x = 0, 1, \dots, n$$

¹Gaussian Hypergeometric Functions