

بِسْمِ اللّٰهِ الرَّحْمٰنِ الرَّحِيْمِ

بِسْمِ اللّٰهِ الرَّحْمٰنِ الرَّحِيْمِ

۳۷۹۷۶



دانشگاه تریست معلم

دانشکده علوم ریاضی و مهندسی کامپیوتر

پایان نامه کارشناسی ارشد

رشته آمار

عنوان:

محاسبه ماکسیمم آنتروپی دیریکلر

برای طرح ریزی داده‌ها

استاد راهنما:

014722

دکتر عین ا... پاشا

نگارش:

نعمت مؤذن

مهر ۱۳۸۰

بسمه تعالی

آگهی دفاع از پایان نامه کارشناسی ارشد آمار

عنوان:

محاسبه ماکسیمم آنتروپی دیریکله برای طرح ریزی داده‌ها

استاد راهنما :	آقای دکتر عین الله پاشا
داور خارجی :	آقای دکتر غلامرضا محتشمی
داور داخلی :	آقای دکتر علی اکبر رحیم زاده
دانشجو :	آقای نعمت مؤذن
زمان :	ساعت ۱۰ صبح روز شنبه مورخ ۸۰/۷/۱۴
مکان :	دانشکده علوم ریاضی و مهندسی کامپیوتر دانشگاه تربیت معلم

خلاصه: اطلاع لگاریتمی که توسط شانون (۱۹۴۸) تعریف شد، کاربردهای زیادی در بخش‌های مختلف آمار پیدا کرده است. از جمله روش تشخیص توزیع نامعلوم مولد داده‌ها می‌باشد.

در این پایان‌نامه ابتدا به تعریف و بررسی بعضی خواص توابع اطلاع از قبیل آنتروپی، چگالی آنتروپی ماکسیمم، اطلاع تمیز، قابلیت تمیز پذیری اطلاع (ID) و اطلاع متناظر می‌پردازیم. سپس مشخصه آنتروپی ماکسیمم از خانواده پارامتری را با طرح پیشین دیریکله برای توزیع نامعلوم مولد داده‌ها ترکیب می‌کنیم که آن را روش آنتروپی ماکسیمم دیریکله (MED) می‌نامیم. روش یافتن توزیع‌های پیشین و پسین MED برای آنتروپی توزیع نامعلوم، پارامترهای گشتاورهای مدل تحت فرض، آنتروپی مدل ME و شاخص ID با استفاده از شبیه‌سازی «مونت کارلو» برای تشخیص توزیع نامعلوم مولد داده‌ها، از مباحثی هستند که در این پایان‌نامه بررسی می‌شوند.

در این پایان‌نامه مقالات زیر مورد استفاده اصلی می‌باشد.

1. Soofi, E.S., N. Ebrahimi and M. Habibullah (1995), "Information Distinguishability with Application to Analysis of Failur Data", Journal of the American Statistical Association, 90, 657-668.
2. Mazzuchi, T. A., Soofi, E.S. and Soyer, R. (2000) "Computation of Maximum Entropy Dirichlet for Modeling Lifetime Data". Computational Statistics and Data Analysis, 32, 361-378.



دانشگاه
علمی

دانشکده علوم ریاضی و مهندسی کامپیوتر

ناریخ
شاره
بیوست
واحد

صورتجلسه دفاع از پایان نامه کارشناسی ارشد

جلسة دفاع از پایان نامه آقای نعمت نعمت مؤذن دانشجوی دوره کارشناسی ارشد رشته ریاضی شاخه آمار تحت عنوان:

محاسبه ماکسیمم آنتروپی دیریکله برای طرح ریزی داده‌ها

در روز شنبه مورخه ۱۴/۷/۸۵ در دانشکده علوم ریاضی و مهندسی کامپیوتر تشکیل گردید و نتیجه آزمون به شرح زیر تعیین می‌گردد. نمره این آزمون ۱۸/۵ (نیم و نیم) می‌باشد.

- | | |
|------------------|-------------------------------------|
| ۱ - عالی | <input checked="" type="checkbox"/> |
| ۲ - بسیار خوب | <input type="checkbox"/> |
| ۳ - خوب | <input type="checkbox"/> |
| ۴ - قابل قبول | <input type="checkbox"/> |
| ۵ - غیرقابل قبول | <input checked="" type="checkbox"/> |

داور داخلی

داور خارجی
دکتر غلامرضا محشمی

استاد راهنمای

دکتر عین‌الله پاکتا

اسماعیل بابلیان
رئیس دانشکده علوم ریاضی و
مهندسی کامپیوتر

تقدیم به :

تقدیم به مادر عزیز و مهربانم
که در تمام مراحل تحصیل مشوق من بوده است.
تقدیم به برادران و خواهران بزرگوارم
که زحمات زیادی را در طول نگارش پایان نامه
به همراه همسران گرامیشان متحمل شدند.

تقدیر و سپاس

لازم است در اینجا از همه کسانی که در تدوین این پایان نامه مرا یاری نمودند سپاسگزاری کنم. از استاد گرانقدر جناب آقای دکتر پاشا که تدوین این پایان نامه با راهنمایی های ایشان صورت گرفت تشکر و قدردانی می نمایم. همچنین از اساتید محترم جناب آقای دکتر رحیم زاده و جناب آقای دکتر محتشمی که قبول زحمت نموده و داوری این رساله را به عهده گرفته‌اند کمال تشکر را دارم. در پایان از دوستان کارشناسی ارشد نیز بخاطر همراهیشان در این تحقیق نهایت قدردانی را می نمایم.

چکیده

اطلاع لگاریتمی که توسط شانون(۱۹۴۸) تعریف شد، کاربردهای زیادی در بخش‌های مختلف آمار پیدا کرده است. از جمله روش تشخیص توزیع نامعلوم مولد داده‌ها می‌باشد.

در این پایان‌نامه ابتدا به تعریف و بررسی بعضی خواص توابع اطلاع از قبیل آنتروپی، چگالی آنتروپی ماکسیمم، اطلاع تمیز، قابلیت تمیز، پذیری اطلاع (ID) و اطلاع متقابل می‌پردازیم. سپس مشخصه آنتروپی ماکسیمم از خانواده پارامتری را با طرح پیشین دیریکله برای توزیع نامعلوم مولد داده‌ها ترکیب می‌کنیم که آن را روش آنتروپی ماکسیمم دیریکله (MED) می‌نامیم. روش یافتن توزیع‌های پیشین و پسین MED برای آنتروپی توزیع نامعلوم، پارامترهای گشتاوری مدل تحت فرض، آنتروپی مدل ME و شاخص ID با استفاده از شبیه‌سازی «مونت کارلو» برای تشخیص توزیع نامعلوم مولد داده‌ها، از مباحثی هستند که در این پایان‌نامه بررسی می‌شوند.

در این پایان‌نامه مقالات زیر مورد استفاده اصلی می‌باشد.

1. Soofi, E.S. , N. Ebrahimi and M. Habibullah (1995), "Information Distinguishability with Application to Analysis of Failure Data", Journal of the American Statistical Association , 90, 657-668.
2. Mazzuchi, T. A, Soofi, E. S. and Soyer, R.(2000) " Computation of Maximum Entropy Dirichlet for Modeling Lifetime Data", Computational Statistics and Data Analysis, 32, 361-378.

فهرست

صفحه	عنوان
۱	پیشگفتار ح
۵	مقدمه و تاریخچه ۱
۵	فصل اول: تعاریف و مقدمات آنتروپی ۵
۱۰	(۱.۱) آنتروپی طرحهای متناهی، عدم حتمیت ۵
۱۰	(۲.۱) ارتباط آنتروپی و اطلاع ۱۰
۱۲	۱.۱. شرح ارتباط آنتروپی و اطلاع ۱۰
۱۷	۲.۱. قضیه یکتایی ۱۲
۱۹	۳.۱. آنتروپی زنجیر مارکف ۱۷
۱۹	(۳.۱) آنتروپی شرطی و اطلاع متقابل ۱۹
۲۱	۱.۱. آنتروپی شرطی ۱۹
۲۴	۲.۱. نتیجه های مهم ۲۱
۲۵	۳.۱. اطلاع متقابل ۲۴
۲۷	۴.۱. چند تعمیم ۲۵
۲۸	۵.۱. جمع بندی مفاهیم ۲۷
۲۸	(۴.۱) آنتروپی متغیر های تصادفی ۲۸
۳۰	۱.۱. تعاریف ۲۸
۳۰	۲.۱. آنتروپی به عنوان اميد ریاضی ۳۰
۳۱	۳.۱. آنتروپی توأم ۳۱

۴.۴.۱. آنتروپی شرطی	۳۱
۵.۴.۱. اطلاع متقابل	۳۱
۶.۴.۱. برخی خواص آنتروپی	۳۲
۷.۴.۱. چند تعمیم	۳۴
۸.۴.۱. آنتروپی تبدیل متغیرهای تصادفی	۳۵
فصل دوم: روش آنتروپی ماکسیمم و توابع توزیع دیریکله	۳۸
(۱.۲) معرفی روش آنتروپی ماکسیمم	۳۸
۱.۱.۲. مقدمه	۳۸
۲.۱.۲. بیان روش	۳۹
۳.۱.۲. قیدهایی به صورت امید ریاضی	۴۰
۳.۱.۲. الف - متغیرهای تصادفی پیوسته	۴۰
۳.۱.۲. ب - متغیرهای تصادفی گستته	۴۲
(۲.۲) طرح دیریکله	۴۴
۱.۲.۲. مقدمه	۴۴
۲.۲.۲. توزیع دیریکله	۴۴
۳.۲.۲. کاربردها	۴۶
۳.۲.۲. الف - برآورد یکتابع توزیع	۴۷
فصل سوم: توابع اطلاع تمیز و خواص آنها	۵۰
(۱.۲) مقدمه	۵۰
(۲.۳) تمیز پذیری اطلاع	۵۳
(۳.۳) تشخیص آماره ها از روی خانواده های پارامتری	۶۰

۱.۳.۳. مقدمه ۶۰	۱.۳.۳
۲.۳.۳. قابلیت تمیز پذیری ID نسبت به نرمال بودن ۶۰	۲.۳.۳
۲.۳.۳. قابلیت تمیز پذیری ID نسبت به نمایی بودن ۶۲	۲.۳.۳
۴.۳) آماره ID ۶۳	(۴.۳)
۵.۳) اطلاع تمیز و اطلاع متقابل ۶۷	(۵.۳)
فصل چهارم: محاسبه ماکسیمم آنتروپی دیریکله برای طرح ریزی داده‌ها ۷۰	
۱.۴) مقدمه ۷۰	(۱.۴)
۲.۴) پیشین MED و آنتروپی کمی شده ۷۳	(۲.۴)
۳.۴) الگوریتم محاسبه ۷۶	(۳.۴)
۴.۴) معرفی روش مونت کارلو ۸۴	(۴.۴)
۵.۴) جمع‌بندی نتایج ۸۶	(۵.۴)
ضمایم ۸۸	
واژه نامه فارسی - انگلیسی ۸۹	
واژه نامه انگلیسی - فارسی ۹۴	
اسامی اشخاص ۹۸	
کتابنامه ۹۹	

پیشگفتار

هدف این تحقیق ارائه روشی بر پایه آنتروپی برای استنتاج بسیزی در مورد اینکه یک توزیع پارامتری به عنوان یک مدل برای داده‌ها مناسب است یا خیر، می‌باشد. در زیر مروری اجمالی بر مباحث بررسی شده در این پایان نامه خواهیم داشت.

فصل اول: تعاریف و مقدمات آنتروپی

در این فصل ابتدا در بخش اول طرح متناهی و عدم حتمیت آن را تعریف کرده و سپس ارتباط میان آنتروپی و اطلاع را در بخش دوم شرح می‌دهیم. با استفاده از یک قضیه مهم اندازه‌ای برای عدم حتمیت به دست می‌آوریم این اندازه در واقع همان اندازه‌ای است که شانون (۱۹۴۸) نخستین بار معرفی نمود. در بخش سوم آنتروپی شرطی و اطلاع متقابل روی طرحها یا فضاهای متناهی بحث می‌شود. بخش چهارم نیز به آنتروپی متغیرهای تصادفی گستته و پیوسته اختصاص یافته است که در آن با استفاده از قضایا، برخی خواص آنتروپی مورد بحث و بررسی قرار گرفته است.

فصل دوم: روش آنتروپی ماکسیمم و توابع توزیع دیریکله

در بخش نخست این فصل ابتدا به معرفی روش آنتروپی‌ماکسیمم می‌پردازیم. سپس قضیه تابع چگالی ماکسیمم آنتروپی را با قیود گشتاوری داده شده بیان و اثبات می‌نماییم. در انتها جدول توابع چگالی آنتروپی‌ماکسیمم با قیود معلوم ارائه می‌شود. در بخش دوم به معرفی توابع توزیع دیریکله که در این پایان نامه استفاده شده است می‌پردازیم آنگاه به کاربرد آن در تعیین برآورد توزیع پیشین برای فرم‌های نامعلوم توابع توزیع خواهیم پرداخت.

فصل سوم: توابع اطلاع تمیز و خواص آنها

این فصل با هدف اندازه گیری تفاوت‌های کلی توزیع‌ها، روی تابع اطلاع «کالبک-لیبلر» بنا شده است. در بخش نخست به توضیح مقدمات استفاده از چگالی آنتروپی ماکسیمم در توابع اطلاع توسط محققان مختلف می‌پردازیم. در بخش بعدی برای مقایسه توزیع‌ها در Ω_θ ، اختلاف اطلاع بین دو چگالی f (نامعلوم) و $*f$ (آنتروپی ماکسیمم) مورد بررسی قرار گرفته است. تابع اطلاع کالبک - لیبلر به عنوان معروف‌ترین اندازه اطلاع نظری اختلاف بین توزیع‌ها شرح داده شده است. نامنفی بودن، جمع پذیری برای متغیرهای تصادفی مستقل از خواص این تابع می‌باشد. سپس شاخص ID توزیع‌ها، با استفاده از تابع اطلاع «کالبک - لیبلر» تعریف می‌شود که قابلیت تبعیض پذیری اطلاع را بیان می‌نماید و یک اندازه نرمال شده می‌باشد. به دست آوردن شاخص ID متناسب با حالتهای نرمال و نمایی از دیگر قسمتهای این بخش است. در بخش چهارم با معرفی آماره ID، برآورد MLE این شاخص ارائه می‌شود. در بخش آخر نیز ارتباط بین اطلاع تمیز و اطلاع متقابل بیان می‌شود.

فصل چهارم: محاسبه ماکسیمم آنتروپی دیریکله برای طرح ریزی داده‌ها

در ابتدا کلاس گشتاوری مورد نظر و چگالی $*f$ در این کلاس معرفی شده سپس توابع اطلاع مورد نظر تعریف می‌شوند. در بخش دوم آنتروپی کمی شده یک توزیع پیوسه f با افزار نمودن خط حقیقی تعریف می‌شود. افزار مفیدی با استفاده از سنگفرش‌های دیریکله ساخته می‌شود. برای تعیین توزیع F نامعلوم از یک طرح پیشین دیریکله استفاده می‌نماییم که یک توزیع کمی شده برای F به دست می‌آید. با استفاده از این توزیع، برآورده برای آنتروپی محاسبه می‌شود. در بخش سوم، آنتروپی کمی شده، پارامترهای گشتاوری، مدل پارامتری و شاخص اطلاع ID با استفاده از پیشین دیریکله، برآورد شده و توزیع پیشین و پسین آنها با استفاده از یک الگوریتم شبیه‌سازی «مونت کارلو» به دست می‌آید.

مقدمه و قاریچه

نظریه اطلاع برای اولین بار در سالهای ۱۹۴۷-۱۹۴۸ توسط "کلود شانون" مطرح شد. این نظریه در آغاز با بسیاری از مسائل ریاضی دشوار رویرو بود. طبیعی بود که شانون و اولین قوانینش که تمام هدف آنها محاسبه نتایج سودمند بود، نمی‌توانست دقت کافی برای این مسائل مشکل ریاضی ارائه کند.

در نتیجه در بسیاری از این تحقیقات، دانشمندان به منظور سهولت اثباتها، مجبور می‌شدند با استدلال از یک ماهیت غیر طبیعی یا مجموعه‌ای از موضوعات به طور مصنوعی (مانند، منابع، کانالها، کدها و غیره) استفاده کنند. بخش عمدۀ ادبیات اولین سالهای نظریه اطلاع بنâچار به محاسبات ناقص ریاضی که مخصوصاً آنرا برای ریاضی دانان بی‌نهایت مشکل می‌کرد، اختصاص یافت.

بعد یک کتاب درسی از نظریه اطلاع توسط "اس. گلدمان" منتشر شد که مثالهای رایجی را در آن بکار برد. سپس مقالات زیادی نوشته شد که پایه‌های ریاضی و محاسباتی نظریه اطلاع را بنا نهادند. به عنوان اولین کارها باید از کار "مک میلان" نام برد که نظریه مفهومات اساسی منابع گسته (منبع، کanal، کدوغیره) را به عنوان اولین تعاریف دقیق ریاضی بکار برد. مهمترین نتیجه این کار برای اثبات قضیه مشهور، باید بکار گرفته می‌شد که در آن هر منبع آرگو دیک گسته، ویژگی را که شانون به منابع نوع مارکف نسبت می‌داد دارا بود که این تقریباً همه محاسبات مجانية نظریه اطلاع را را لایه می‌داد. این وضع به تمام قسمتهای نظریه اطلاع گسته اجازه می‌داد که بدون محدودیت، همانند نظریه شانون، به منابع نوع مارکف ساخته شوند. "مک میلان" در ادامه مقاله‌اش سعی نمود قضیه اساسی شانون را روی کانالهای بانوفه روی اصل دقیقی پایه ریزی کند که در انجام آن معلوم شد که اثبات خلاصه وار که توسط

شانون داده شده بود دارای وقه و نقص زیادی است که حتی در حالت منابع مارکف هنوز باقی مانده بود.

سپس لازم است که به کار "فین اشتین" توجه شود که او مثل "مک میلان" فرض کرد که قضیه شانون روی کانالهای بانوفه، متنهای نظریه کلی اطلاع گستته باشد و متعهد به دادن اثباتی دقیق از این نظریه شد.

با پذیرفتن عملیات ریاضی "مک میلان" با اجتناب از پیروی مسیر اصلی شانون، اثباتی را بنا نهاد. که در آن از یک عقیده، نتیجه بخش کاملاً جدید آشکار تحت عنوان "مجموعه قابل تمیز از دنباله‌ها" استفاده نمود.

هر چند "فین اشتین" اثبات آنرا برای مسائل و مسائلی که دارای بار عملی کمتری بودند بکار برد، که در آن سیگنانالهای پی در پی منبع دو بدو مستقل هستند و حافظه کانال نیز صفر می‌باشد. در حالت کلی تر، فقط به طور خلاصه نشان داد که چطور می‌توان استدلال لازم را برای مستقل بودن انجام داد. متأسفانه از نحاط معنایی مشکلات زیادی باقی ماند.

کالبک ولیلر (۱۹۵۱) اندازه فاصله بین دو چگالی مختلف را تعریف نمودند که نزدیکی و اطلاع تمیز دو چگالی را نشان می‌داد. این تابع آنتروپی یک چگالی نسبت به چگالی دیگر نیز می‌باشد.

جینس (۱۹۵۷) اصل آنتروپی ماکسیمم را بیان نمود. این اصل چگالی آنتروپی ماکسیمم را با قیود گشتاوری داده شده مشخص می‌کرد و به کاربردهای زیادی در علوم و زمینه‌های مختلف از جمله مکانیک آماری، نجوم، هواشناسی، اقتصاد، جغرافیا، تجارت بین المللی، بانکداری، توزیع جمعیتی، پژوهشکی و غیره دست یافت.

کالبک (۱۹۵۹) تابع اطلاع تمیز بین دو چگالی را با استفاده از اصل آنتروپی ماکسیمم بیان نمود که نزدیکی چگالی نامعلوم و چگالی آنتروپی ماکسیمم را نشان می‌دهد این پایان نامه براساس تابع اطلاع کالبک لیلر (۱۹۵۹) بنا شده است و هدف آن تشریح توابع اطلاع و ارائه

روشی برای برآورد آنتروپی، توابع اطلاع، پارامترهای مدل و توزیعهای پیشین و پسین آنها با بکارگیری «ماکسیمم آنتروپی دیریکله» می‌باشد.