

کاربرد آمار در داده کاوی

مقدمه و مقایسه

آمار شاخه ای از علم ریاضی است که به جمع آوری توضیح و تفسیر داده ها می پردازد. [۳] این مبحث به گونه ای است که روزانه کاربرد زیادی دارد. در مقایسه این علم با data mining قدمت بیشتری دارد و جزء ورشهای کلاسیک داده کاوی محسوب می شود. وجه اشتراک تکنیکهای آماری و data mining بیشتر در تخمین و پیش بینی است. [۲] البته از آزمونهای آماری در ارزیابی نتایج داده کاوی نیز استفاده می شود. در کل اگر تخمین و پیش بینی جزء وظایف data mining در نظر گرفته شوند، تحلیل های آماری، data mining را بیش از یک قرن اجرا کرده است. به عقیده بعضی DM ابتدا از آمار و تحلیل های آماری تحلیل شروع شد. [۲] می توان تحلیل های آماری از قبیل فاصله اطمینان، رگرسیون و... را مقدمه و پیش زمینه DM را دانست که بتدریج در زمینه های دیگر و متد های دیگر رشد و توسعه پیدا کرد. پس در واقع متدهای آماری جزو روشهای کلاسیک و قدیمی DM محسوب می شوند. در جایی اینگونه بحث می شود که با تعریف دقیق، آمار یا تکنیکهای آماری جزء داده کاوی (data mining) نیستند. این روشها خیلی قبل تر از data mining استفاده می شدند. با این وجود، تکنیکهای آماری توسط داده ها بکار برده می شوند و برای کشف موضوعات و ساختن مدل های پیشگویانه مورد استفاده قرار می گیرند. [۳]

همانگونه که واضح و مشخص است با گذشت زمان علم نیز پیشرفت می کند، هر چه به جلوتر می رویم روشهای جدید تر و بهتر مورد استفاده قرار می گیرد، علم امروز نسبت به دیروز جدیدتر است. روشهای جدید علمی در پی کشف محدودیتهای روشهای قدیمی ایجاد می شود، و از آنجایی که روشهای آماری جزء روشهای قدیمی Data mining محسوب می شوند، از این قاعده کلی که دارای محدودیت هستند مستثنی نیستند. داشتن فرض اولیه در مورد داده ها، یکی از این موارد است. در اینجا به تشریح بیشتر تفاوت های بین مباحث و متدهای آماری و دیگر متدهای داده کاوی که در کتابهای مختلف بحث شده است می پردازیم.

تکنیکهای داده کاوی و تکنیکهای آماری در مباحثی چون تعریف مقدار هدف برای پیش گویی، ارزیابی خوب و داده های دقیق (تمیز) (clean data) خوب عمل می کنند، همچنین این موارد در جاهای یکسان برای انواع یکسان از مسایل (پیش گویی، کلاس بندی و کشف) استفاده می شوند، بنابراین تفاوت این دو چیست؟ چرا ما آنچنان که علاقه مند بکاربردن روشهای داده کاوی هستیم علاقه مند

روشهای آماری نیستیم؟ برای جواب این سوال چندین دلیل وجود دارد اول اینکه روشهای کلاسیک داده کاوی از قبیل شبکه های عصبی، تکنیک نزدیک ترین همسایه روشهای قوی تری برای داده های واقعی به ما می دهند و همچنین استفاده از آنها برای کاربرانی که تجربه کمتری دارند راحت تر است و بهتر می توانند از آن استفاده کنند. دلیل دیگر اینکه بخاطر اینکه معمولاً داده ها اطلاعات زیادی در اختیار ما نمی گذارند، این روشها با اطلاعات کمتر بهتر می توانند کار کنند و همچنین اینکه برای داده ها وسیع کاربرد دارند. [۳]

در جایی دیگر اینگونه بیان شده که داده های جمع آوری شده نوعاً خیلی از فرضهای قدیمی آماری را در نظر نمی گیرند، از قبیل اینکه مشخصه ها باید مستقل باشند، تعیین توزیع داده ها، داشتن کمترین همپوشانی در فضا و زمان اغلب داده ها هم پوشانی زیاد می دارند، تخلف کردن از هر کدام از فرضها می توان مشکلات بزرگی ایجاد کند، زمانی که یک کاربر (تصمیم گیرنده) سعی می کند که نتیجه ای را بدست آورد. داده های جمع آوری شده بطور کلی تنها مجموعه ای از مشاهدات چندی بعد است بدون توجه به اینکه چگونه جمع آوری شده اند [۵].

در جایی پایه و اساس Data mining به دو مقوله آمار و هوش مصنوعی تقسیم شده است که روشهای مصنوعی به عنوان روشهای یادگیری ماشین در نظر گرفته می شوند. فرق اساسی بین روشهای آماری و روشهای یادگیری ماشین (machine learning) بر اساس فرضها و یا طبیعت داده هایی که پردازش می شوند. بعنوان یک قانون کلی فرضها تکنیکهای آماری بر این اساس است که توزیع داده ها مشخص است که بیشتر موارد فرض بر این است که توزیع نرمال است و در نهایت درستی یا نادرستی نتایج نهایی به درست بودن فرض اولیه وابسته است. در مقابل روشهای یادگیری یادگیری ماشین از هیچ فرض در مورد داده ها استفاده نمی کند و همین مورد باعث تفاوتی بین این دو روش می شود.

به هر حال ذکر این نکته ضروری به نظر می رسد که بسیاری از روشهای یادگیری ماشین برای ساخت مدل dataset از حداقل چند استنتاج آماری استفاده می کنند که این مساله بطور خاص در شبکه عصبی دیده می شود. [۱]

بطور کلی روشهای آماری روش های قدیمی تری هستند که به حالت های احتمالی مربوط می شوند. Data mining جایگاه جدید تری دارد که به هوش مصنوعی یادگیری ماشین سیستمهای اطلاعات مدیریت (MIS) و متدلوژی Database مربوط می شود.

روشهای آماری بیشتر زمانی که تعداد دادهها کمتر است و اطلاعات بیشتری در مورد داده ها می توان بدست آورد استفاده می شوند به عبارت دیگر این روشها با مجموعه داده های کوچک تر سر و کار دارند همچنین به کاربران ابزارهای بیشتری برای امتحان کردن داده ها با دقت بیشتر فهمیدن ارتباطات بین داده ها می دهد. بر خلاف روشهایی از قبیل شبکه عصبی که فرآیند مبهمی دارد. پس به طور کلی این روش در محدوده مشخصی از داده های ورودی بکار می رود. بکار بردن این روشها مجموعه داده های مجموعه داده های زیاد احتمال خطا در این روشها را زیاد می کند. چون در داده ها احتمال noise و خطا بیشتر می شود و نیز روشهای آماری معمولاً به حذف noise می پردازند، بنابراین خطای محاسبات در این حالت زیاد می شود. [۷]، [۸]

در بعضی از روشهای آماری نیاز داریم که توزیع داده ها را بدانیم. اگر بتوان به آن دسترسی پیدا کرده با بکار بردن روش آماری می توان به نتایج خوبی رسید.

روشهای آماری چون پایه ریاضی دارند نتایج دقیق تری نسبت به دیگر روشهای Data mining آریه می دهند ولی استفاده از روابط ریاضی نیازمند داشتن اطلاعات بیشتری در مورد داده ها است.

مزیت دیگر روشهای آماری در تعبیر و تفسیر داده ها است. هر چند روشهای آماری به خاطر داشتن ساختار ریاضی تفسیر سخت تری دارند ولی دقت نتیجه گیری و تعبیر خروجی ها در این روش بهتر است بطور کلی روشهای آماری زمانی که تفسیر داده ها توسط روشهای دیگر مشکل است بسیار مفید هستند.

تفاوتهای کلی روشهای آماری و دیگر روشهای Data mining در جدول آریه شده است:

روشهای آماری	دیگر روشهای Data mining
داشتن فرض اولیه	بدون فرض اولیه

تنها برای داده های عددی کاربرد دارند	در انواع مختلفی از داده ها کاربرد دارند نه فقط داده های عددی
در محدوده کوچکی از داده ها	در محدوده وسیع تری از داده ها
حذف noise ها ، داده های نامشخص و فیلتر کردن dirty data	Data mining به داده های درست clean data بستگی دارند
روشهای رگرسیون و استفاده از معادلات	استفاده از شبکه عصبی
استفاده از چارتهای دو بعدی و سه بعدی	استفاده از Data visualization
استفاده از روابط ریاضی	استفاده از روشهای یادگیری ماشین و هوش مصنوعی
در cluster analysis و descriptive statistical کاربرد دارد.	در یادگیری غیر نظارتی کاربرد بیشتر دارد

همچنین می توان گفت که در DM داده ها اغلب بر اساس همپوشانی نمونه هاست، نسبت به اینکه بر اساس احتمال داده ها باشد. همپوشانی نمونه ها برای آشنایی همه انواع پایه ها برای تخمین پا را مترها مشهور است. و همچنین اغلب استنتاج های آماری نتایج ممکن است مشارکتی باشد تا اینکه سببی باشند. تکنیکهای ماشین را به سادگی می توان تفسیر کرد. مثلاً روش شبکه عصبی بر اساس یک مدل ساده بر اساس مغز انسان استوار است. یعنی همان ساختار مغز انسان را اجرا می کنند ولی خروجی های بسیاری از روشهای آماری ساختار ریاضی دارند، مثلاً یک معادله است که تعبیر و تفسیر آن مشکل تر است. در مورد روش های آماری باید این مطلب را گفت بدون توجه به اینکه مدل کاربردی، مدل آماری است یا خیر، تستهای آماری می تواند برای تحلیل نتایج مفید باشد.

با ارایه توضیحات داده شده درباره های تفاوت های روشهای آماری و دیگر روشهای DM در ادامه به کاربردهای روش روشهای آماری و بحثهای مشترک آنها رو DM می پردازیم .

کاربردهای روشهای آماری:

Data mining معمولاً وظایف یا به عبارت بهتر استراتژیهای زیر را در داده‌ها بکار می‌برد:

- توضیح و تفسیر (description)
- تخمین (estimation)
- پیش‌بینی (prediction)
- کلاس‌بندی (classification)
- خوشه‌سازی (clustering)
- وابسته‌سازی و ایجاد رابطه (association)

در جدول زیر استراتژی‌ها و روشهای هر استراتژی مشخص شده است:

روشها	استراتژیها
تحلیل داده‌ها	توضیح و تفسیر
تحلیل‌های آماری	تخمین
تحلیل‌های آماری	پیش‌بینی
الگوریتم نزدیک‌ترین همسایه	کلاس‌بندی
درخت تصمیم	کلاس‌بندی
شبکه‌های عصبی	کلاس‌بندی
خوشه‌سازی k-mean	خوشه‌سازی
شبکه‌های kohonen	خوشه‌سازی
وابسته‌سازی و ایجاد رابطه	رابطه‌سازی

البته باید گفت که روشهای data mining تنها به یک استراتژی خاص محدود نمی‌شوند و نتایج یک را همپوشانی بین روشها نشان می‌دهد. برای مثال درخت تصمیم ممکن است که در کلاس‌بندی

تخمین و پیش بینی کاربرد داشته باشد. بنابراین این جدول را نباید به عنوان تعریف تعریف تقسیم بندی از وظایف در نظر گرفته شود بلکه به عنوان یک خروجی از آنچه که ما به عنوان وظایف datamining آشنایی پیدا کردیم در نظر گرفته می شود.

همانگونه که از جدول پیداست روشهای آماری در مباحث تخمین و پیش بینی کاربرد دارند. در تحلیل آماری تخمین و پیش بینی عناصری از استنباطهای آماری هستند. استنباطهای آماری شامل روشهایی برای تخمین و تست فرضیات درباره جمعیتی از ویژگیها براساس اطلاعات حاصل از نمونه است. یک جمعیت شامل مجموعه ای از عناصر از قبیل افراد ایتیم یاداده ها می که دریک مطالعه خاص آمده است. بنابراین در اینجا به توضیح این دو استراتژی می پردازیم.

۱- تخمین:

در تخمین به دنبال این هستیم که مقدار یک مشخصه خروجی مجهول را تعیین کنیم، مشخصه خروجی در مسایل تخمین بیشتر عددی هستند تا قیاسی [۱]. بنابراین مواردی که بصورت قیاسی هستند باید به حالت عددی تبدیل شوند. مثلاً موارد بلی، خیر به 0 و 1 تبدیل می شود.

تکنیکهای نظارتی DM قادرند یکی از دو نوع مسایل کلاس بندی یا تخمین را حل کنند، نه اینکه هر دو را. یعنی اینکه تکنیکی که کار تخمین را انجام می دهد، کلاس بندی نمی کند.

روشهای آماری مورد استفاده در این مورد بطور کلی شامل تخمین نقطه و فاصله اطمینان میباشد. تحلیل های آماری تخمین و تحلیل های یک متغیره و... از این جمله می باشند.

در توضیح اینکه چرا به سراغ تخمین می رویم باید گفت که مقدار واقعی پارامترها برای ما ناشناخته است. مثلاً مقدار واقعی میانگین یک جامعه مشخص نیست. داده ها ممکن است که بطور رضایت بخشی جمع آوری نشده باشد یا به عبارتی warehouse نشده باشد. به همین دلیل تحلیل گران از تخمین استفاده میکنند.

در خیلی از موارد تعیین میانگین مجموعه ای از داده ها برای ما مهم است. مثلاً میانگین نمرات درسی یک کلاس، میانگین تعداد نفراتی که در یک روز به بانک مراجعه می کنند، متوسط مقدار پولی که افراد دریک شعبه خاص از بانک واریز می کنند و موارد اینچنینی.

زمانی که مقدار یک آماره را برای برآورد کردن پارامتر یک جامعه به کار ببریم، آن پارامتر را تخمین زده ایم، و به مقدار این آماره برآورد نقطه ای پارامتر اطلاق می کنیم. در واقع از کلمه نقطه برای تمایز بین برآورد کننده های نقطه ای و فاصله ای استفاده می کنیم. از مهمترین تخمین زننده ها است که به ترتیب برآورد واریانس و میانگین جامعه هستند. خود برآورد کننده ها دارای خاصیت هایی چون ناریبی، کارایی، ناسازگاری، بسندگی و... هستند، که هر یک به بیان ویژگی خاصی از آنها می پردازند و میزان توانایی آنها را در تخمین درست و دقیق یک پارامتر تعیین می کنند.

در تخمین نیازمند داشتن اندازه نمونه هستیم، در تعیین اندازه نمونه می توان از رابطه زیر استفاده کرد:

که p احتمال رخداد و e درصد خطای پذیرفته شده است که در اینجا 5% در نظر می گیریم.

پر کاربرد ترین تخمین زننده، تخمین زننده میانگین جامعه است، ساده ترین رابطه ای که برای میانگین داده ها می توان نوشت بدین صورت است:

که n تعداد نمونه ها و مقدار هر نمونه است. در اینجا تمام نمونه ها ارزش یکسانی دارند ولی گاهی اوقات نیاز است که نمونه ها بر اساس اهمیتی که دارند وزن دهی شوند.

بدین صورت :

که ها در اینجا وون هر نمونه A ام هستند. در این حالت برای تعیین مجموع اوزان نمونه ها به جای n ، می باشد.

در مواردی نیز تخمین فاصله برای ما اهمیت دارد. فاصله اطمینان شامل فاصله ای است که با درصدی از اطمینان می توانیم بگوییم که مقدار یک پارامتر درون این فاصله قرار می گیرد. به عبارت دیگر اگر چه برآورد نقطه ای طریقه متداول توصیف برآورد هاست اما درباره آن، جا برای پرسشهای زیادی باقی است. مثلاً برآورد نقطه ای به ما نمی گوید که برآورد بر چه مقداری از اطلاعات مبتنی است. و چیزی درباره خطا بیان نمی کند. بنابراین می توانیم که برآورد پارامتر را با بعلاوه کردن اندازه کردن اندازه نمونه و مقدار واریانس، یا اطلاعات دیگری درباره توزیع نمونه گیری کامل کنیم. این کار ما را قادر می سازد که اندازه ممکن خطا را برآورد کنیم.

یک برآورد فاصله ای، فاصله ای به شکل است که در آن و مقادیر متغیرهای تصادفی مناسبی برای هستند، منظور از مناسب آن است که به ازای احتمال مشخصی مانند داریم:

برای مقدار مشخص ، را یک فاصله اطمینان برای می نامیم. همچنین ، درجه اطمینان ، و دو سر فاصله کرانهای اطمینان پایینی و بالایی نامیده می شود. مثلاً برای ، درجه اطمینان 95% است و یک فاصله اطمینان 95% بدست می آوریم. فاصله اطمینان از اکثر توزیع ها ، همانند توزیع نرمال ، خی دو، t استودنت و توزیع F و ... استفاده می کند. مثلاً اگر مقدار میانگین یک نمونه تصادفی به اندازه n از جامعه نرمال و وایانس معلوم باشد آنگاه

یک فاصله اطمینان برای میانگین جامعه است. [۶]

در خیلی از موارد تعیین نقطه دقیق یک پارامتر ممکن نیست، ولی فاصله اطمینان ، اطمینان ما را از قرار گرفتن مقدار پارامتر در یک بازه تضمین می کند. فاصله اطمینان را می توان برای اکثر توزیع ها از جمله توزیع خی دو، توزیع t استودنت و توزیع F و ... بدست آورد.

۲- پیش بینی (prediction) :

هدف از انجام پیش بینی تعیین ترکیب خروجی با استفاده از رفتار موجود می باشد. یعنی در واقع رسیدن به یک نتیجه بوسیله اطلاعات موجود از داده ها. مشخصه های خروجی در این روش هم می توانند عددی باشند وهم قیاسی. [۱] این استراتژی در بین استراتژی های data mining از اهمیت خاصی برخوردار است، و مفهوم کلی تری را نسبت به موارد دیگر دارد. خیلی از تکنیکهای نظارتی data mining که برای کلاس بندی و تخمین مناسب هستند در واقع کار پیش بینی انجام می دهند. آنچه از کتابهای آماری و data minig تحت عنوان پیش بینی برمی آید رگرسیون و مباحث مربوط به آن است. در واقع در اکثر این کتابها هدف اصلی از انجام تحلیل های آماری برای داده کاوی، رگرسیون داده هاست و این بعنوان وظیفه اصلی متد های آماری معرفی می شود.

اهداف تحلیل رگرسیون:

با انجام رگرسیون می خواهیم اهداف زیر را دنبال کنیم:

۱- بدست آوردن رفتار متغیر Y توسط متغیر X ، یعنی اینکه متغیر Y با تغییر X در نمونه ها چه رفتاری را از خود نشان می دهد. مثلاً در نمونه ای این رفتار خطی است یا اینکه شکل منحنی خواهد داشت.

۲- پیش بینی بر اساس دادهها برای نمونه های آینده، که هدف اصلی در داده کاوی از طریق متدهای آماری است. مثلاً از روی اطلاعاتی مثل داشتن کارت اعتباری یک فرد جدید، نوع جنسیت او، سن فرد، میزان درآمد سالیانه او بتوان حدس زد که این فرد از بیمه عمر استفاده می کند یا خیر. و یا اینکه با داشتن اطلاعات در مورد داشتن یا نداشتن کارت اعتباری و بیمه عمر، سن فرد بتوان جنسیت فرد را تعیین کرد.

۳- استنباط استنتاجی یا تحلیل حساسیت، تعیین اینکه اگر X به اندازه خاصی تغییر کند Y تا چه اندازه تغییر خواهد کرد. هدف از فهمیدن اینکه چگونه تغییرات Y تابعی از X است. باید توجه داشت که نوع تغییرات مدل گرسیونی خاصی را می دهد.

اهداف مدلسازی برای تشریح ارتباط بین X و Y استفاده از نتایج مدل برای پیش بینی کاربردهای تخمین عبارت است. اما استنباط استنتاجی یک مقوله ظریف تری است. زمانی که به استنباط آماری فکر میکنیم در واقع درباره تغییر رفتاری و متغیرهای کنترل فکر می کنیم.

متغیرهای رفتاری مشخصه هایی را ارائه میکنند که تبحر و تجربه خاصی دارند یا اینکه قابلیت آن بحر را دارند. مثلاً مقدار دز دارو که برای بیمار استفاده می شود در تجربه پزشکی. همچنین متغیرهای کنترل دیگر ویژگی ها در یک محیط آزمایشی را اندازه میگیرند، از قبیل وزن بیمار که قبل از رفتار اندازه گیری می شود.

اگر ما برای یکی از متغیرهای رفتاری، کنترل انجام دهیم، رگرسیون ما احتمالاً استنباط های استنتاجی را درست حدس میزند. و اگر ما علاقه مند به هر دو مورد پیش بینی انتخاب سهم و تخمین اثرات علتها باشیم تایید هر دو مورد را بعنوان متغیرهای خروجی که همپوشانی دارند در نظر می گیریم.

روشهای مختلف رگرسیون برای داده کاوی وجود دارد. رگرسیون خطی بیشترین کاربرد را دارد و همچنین مشتقات آن حایز اهمیت است. یک نمونه از آن مشتقات آن رگرسیون خطی سلسله مراتبی یا رگرسیون چند سطحی است. این روش یکی از ابزارهای تحلیل داده‌های پیچیده از قبیل افزایش فرکانس در تحقیقات مقدماتی را شامل می‌شود. مدل‌های رگرسیون چند سطحی برای حالت‌هایی که همپوشانی در سطوح مختلف وجود دارد مفید است. برای مثال اطلاعات آموزشی ممکن است اطلاعاتی از قبیل اطلاعات فردی دانش آموزان (نام، نام خانوادگی و در کل پیش زمینه خانوادگی)، اطلاعات سطح کلاس از قبیل ویژگی‌های معلم و همچنین اطلاعات درباره مدرسه همانند سیاست آموزشی و... باشد. حالت دیگر مدل‌های چند سطحی، تحلیل داده‌های بدست آمده از نمونه‌های خوشه بندی شده است. یک خانواده از مدل‌های رگرسیون، به عنوان متغیرهای شاخص بر رتبه بندی یا خوشه بندی است علاوه بر اینکه همپوشانی را اندازه می‌گیرد. با نمونه خوشه بندی شده مدلسازی چند سطحی برای توسعه نمونه‌هایی که داخل خوشه نیستند، لازم است. [۴]

در روش رگرسیون چند سطحی یا سلسله مراتبی محدودیتی برای تعداد سطوح تغییر که می‌تواند انجام شود، وجود ندارد روش‌های بیزی در تخمین پارامترهای مجهول کمک می‌کند، هر چند که محاسبات پیچیده‌ای دارد. ساده‌ترین توسعه از رگرسیون همپوشانی مجموعه‌ای از متغیرهای شاخص برای کلاس بندی نمونه‌های آموزشی یا رتبه بندی و خوشه بندی در نمونه‌های داده شده است. همچنین به عنوان توسعه رگرسیون خطی در نظر گرفته می‌شود، که در ادامه به توضیح آن می‌پردازیم [۴]

۱- رگرسیون خطی (Linear regression)

یکی از هدف‌های اصلی بسیاری از پژوهش‌های آماری ایجاد وابستگی‌هایی است تا پیش بینی یک یا چند متغیر را بر حسب سایرین ممکن می‌سازد. مثلاً مطالعاتی انجام می‌شود تا فروش‌های بالقوه یک محصول جدید را بر حسب قیمت آن، وزن یک بیمار را بر حسب تعداد هفته‌هایی که پرهیز داشته است، پیش بینی کند.

در عمل مسایل متعددی وجود دارند که در آن‌ها مجموعه‌ای از داده‌ها زوج شده بر آن دلالت می‌کند که رگرسیون خطی است و در آن توزیع توأم متغیرهای تصادفی تحت بررسی رانمی دانیم اما با این حال می‌خواهیم که ضرایب رگرسیون را برآورد کنیم.

روش رگرسیون خطی یک تکنیک یادگیری نظارتی است که به وسیله آن می‌خواهیم تغییرات یک متغیر وابسته بوسیله ترکیب خطی از یک یا چند متغیر مستقل مدل کنیم. حالت کلی معادله آن به این صورت است:

$$(1) \quad f(x_1+x_2+\dots+x_n)=a_1x_1+a_2x_2+\dots+a_nx_n+b$$

که x ها متغیر مستقل و a ها و b ضرایب ثابت هستند و $f(x_1, x_2, \dots, x_n)$ متغیر وابسته می‌باشند. حالت ساده این معادله بصورت $y = ax + b$ (۲) است که در اینجا y متغیر وابسته است به حالت ساده شده معادله ۱ (یعنی معادله ۲) $\text{shope-intercept from}$ می‌گویند.

یک روش برای تعیین ضرایب a, b روش حداقل مربعات است. ملاک کمترین مربعات این است که مجموع مربعات انحراف‌ها را مینیمم کنیم؛ بنابراین اگر مجموعه‌های از داده‌های زوج شده مانند $\{(x_i, y_i), i=1, 2, \dots, n\}$ داده شده باشد، برآوردهای کم‌ترین مربعات ضرایب رگرسیون، مقادیری مانند a, b هستند که به ازای آنها کمیت

مینیمم است:

e در شکل مشخص شده است:

$$e_i = y_i - (ax_i + b)$$

بنابر این در حالت ساده اگر یک نمونه n تایی داشته باشیم مقادیر a, b را از طریق روابط زیر برآورد می کنیم :

مزیت رگرسیون خطی این است که فهمیدن و کار با آن ساده است در حالت کلی برای استراتژی و پیش بینی مناسب است. با بکار بردن این روش از نتایج خروجی می توان دریافت که این روش مناسب بوده یا خیر . بنابر این معیارهایی داریم که با استفاده از آنها می توان دریافت که آیا می توان به نتایج خروجی اطمینان کرد یا خیر.

آنچه در انجام رگرسیون مهم به نظر می رسد، تعیین میزان همبسته بودن داده ها به یکدیگر است. با مشخص کردن میزان همبسته بودن داده های متغیرهای ورودی و خروجی می توان دریافت که رگرسیون خطی برای انجام داده کاوی مناسب است یا خیر، بنابراین ضریب همبستگی و برآوردهای آن در بسیاری از پژوهشهای آماری اهمیت دارند. شرایطی که وقتی چند متغیر پیشگو (X_i) با یکدیگر هم پوشانی دارند، این هم پوشانی منجر ناستواری و تزلزل در فضای جواب می شود، همچنین منجر به نتایج بی ارتباط (بی ربط) می شود. حتی اگر از این تزلزل اجتناب شود هم پوشانی بین متغیرهایی که میزان بین متغیرهایی همبستگی آنها زیاد است، منجر به تاکید کردن روی بخش خاصی از مدل می شود. [۲]

بنابر این از بین متغیر های ورودی مواردی که با هم بستگی زیادی دارند، نباید با هم در تعیین ارزش متغیر خروجی بکار بره شوند. و از طرفی کاربرد رگرسیون خطی منوط به همبستگی متغیر های ورودی و خروجی است. در تحلیل همبستگی نرمال مربوط به داده های زوج شده، با استفاده از روابط ریاضی می توان به، که بیانگر ضریب همبستگی نمونه ای است رسید. رابطه چنین است :

که به ترتیب میانگین متغیر های ورودی و خروجی هستند . را معمولاً با r نمایش می دهند و رابطه ساده شده آن به این صورت است:

که: و شدت بستگی بین X, Y را اندازه می گیرد

در صورتی که $r=0$ باشد، این دو متغیر (X, Y) نسبت به هم نا همبسته اند، و هر چه صفر فاصله می گیرد، بطرف $+1$ و -1 میزان همبسته بودن آنها زیادتر می شود، و $+1$ همبستگی خطی مثبت و -1 همبستگی خطی منفی را نشان می دهد. در حالت توزیع نرمال دو متغیر، صفر بوده $r(=0)$ مستقل بودن این دو متغیر را از هم نشان می دهد. رابطه روبرو را در نظر بگیرید:

وقتی که باشد، نتیجه می شود که و این بدان معنی است که همبستگی خطی کاملی بین X و Y موجود است. برای تفسیر مقادیر r بین 0 و $+1$ یا 0 و -1 ، این معادله را نسبت به حل کرده نتیجه را در 100 ضرب می کنیم بنابراین داریم:

که در آن تغییر کلی Y ها و تغییر شرطی Y ها را به ازای مقادیر ثابت X اندازه می گیرند. بنابراین آن قسمت از کل Y ها که در اثر بستگی به X قابل توضیح است اندازه می گیرد. پس $100r^2$ درصد تغییر کلی از Y ها است که در اثر بستگی به X قابل توضیح است.

مثلاً وقتی $r=5\%$ در این صورت 25% از تغییر Y هاست که در اثر بستگی به X قابل توضیح است. و وقتی $r=7\%$ در این صورت 49% درصد از تغییر Y ها در اثر بستگی به X قابل توضیح است. بنابراین می توانیم بگوییم که یک همبستگی $r=7\%$ تقریباً دو برابر قوی تر یک همبستگی $r=5\%$ است. [۱]

همچنین تحلیل رگرسیونی نرمال برای حالت چند گانه رابطه مفید زیر را که بر اساس توزیع t بدست می دهد، ارائه می کند:

که در این رابطه ضریب متغیر X_i و مقدار عدد ثابت در معادله رگرسیون است.

n تعداد نمونه ها، k تعداد متغیرهای ورودی (X_i ها) است.

همچنین در نظر بگیرید که یک سری داده متشکل از k متغیر ورودی و یک متغیر خروجی که تعداد هر نمونه از متغیر n تا باشد، با ضرایب $i=0, 1, \dots, k$ بصورت ماتریس به شکل زیر نمایش داده شود:

همچنین را ترانهاده و را نهاده (و) $B=$ در نظر بگیرید. بدین ترتیب و C_{ii} درایه ماتریس معکوس X یعنی است. باید درایه C_{22} در این ماتریس را حساب کنیم.

عبارت t دارای توزیع t با $n-k-1$ درجه آزادی است. که یک آماره مناسب برای آزمون میزان تأثیری که ضریب هر t یعنی در معادله رگرسیونی دارد. [۶]

۲- Logistic Regression

این روش یکی از تکنیکهای یادگیری نظارتی و در حالتی که نتایج خروجی به صورت binary هستند، مورد توجه قرار می گیرد. در کل زمانی نتایج خروجی به صورت binary هستند رگرسیون خطی خیلی کارا نیست، در این حالت استفاده از این تکنیک مناسب تر است. نکته دیگر اینکه این روش یک تکنیک رگرسیون غیر خطی است و لزومی ندارد که داده ها حالت خطی داشته باشند. اگر بخواهیم دلیل استفاده Logistic regression را بیان کنیم باید اینگونه بحث کنیم در رگرسیون خطی علاوه بر اینکه نتایج خروجی باید به صورت عددی باشد، متغیرها هم باید به صورت عددی باشد بنابراین حالتی که به صورت کتگوری (قیاسی) هستند باید به حالت عددی تغییر شکل پیدا کنند. مثلاً جنسیت افراد از حالت زن و مرد بوده به ترتیب به 0 و 1 تغییر پیدا می کند. در این روش اگر نتایج خروجی (متغیر خروجی) بصورت binary باشد می تواند مفید باشد. چون اساس رگرسیون خطی در این حالت ایراد پیدا می کند و ارزش قیدی که بر روی متغیر وابسته قرار می گیرد توسط معادله رگرسیون در نظر گرفته نمی شود.

در واقع چون رگرسیون خطی، معادله یک خط را ترسیم می کند، نمی تواند حالت مثبت و منفی یا به عبارتی صفر و یک را در نظر بگیرد.

به همین دلیل برای اینکه بتوان حالت‌های binary را هم در نظر گرفت، باید شکل معادله را تغییر داد. با این تغییر شکل معادله رگرسیون احتمال اتفاق افتادن یا اتفاق نیفتادن یک واقعه را بدست می‌دهد.

با تغییر شکل رگرسیون خطی به حالت Logistic regression این مشکل حل می‌شود.

معادله خطی را می‌توان بدین صورت نوشت:

که بیانگر احتمال اینکه متغیر وابسته (y) مقدار 1 را بگیرد به شرط اینکه ترکیبی از X را داشته باشیم. بصورت کلی تر و برای حذف محدودیت‌های مسئله‌ها حالت احتمالی $y=1$ را نسبت به $y=0$ در

نظر می‌گیریم یعنی به صورت

ولگاریتم طبیعی این عبارت را برابر قرار می‌دهیم که X برداری بصورت است، و در نهایت از رابطه \ln

بالا مقدار بدست می‌آید که برابر است با:

این معادله، معادله Logistic regression را تشکیل می‌دهد.

اگر بخواهیم منحنی این معادله را نشان دهیم بصورت روبرو می‌باشد:

که بیانگر غیر خطی بودن این معادله است. در نهایت برای فهم بهتر مسأله مثالی ارائه می‌شود.

داده‌های زیر را وارد Excel کرده و ضرایب متغیرها و مقدار ثابت b را توسط تابع LINEST بدست

می‌آوریم داده‌ها و نتایج به این صورت می‌باشد:

instance	income	Credit card insurance	sex	age	Life insurance promotion	Computed probability
1	40	0	1	45	0	0.007
2	30	0	0	40	1	0.987
3	40	0	1	42	0	0.024
4	30	1	1	43	1	1.000
5	50	0	0	38	1	0.999
6	20	0	0	55	0	0.049
7	30	1	1	35	1	1.000
8	20	0	1	27	0	0.584

9	30	0	1	43	0	0.005
10	30	0	0	41	1	0.981
11	40	0	0	43	1	0.985
12	20	0	1	29	1	0.380
13	50	1	0	39	1	0.999

این مثال ۴ مشخصه ورودی و یک مشخصه خروجی دارد که ضرایب متغیرهای ورودی در زیر محاسبه شده است:

$$ax+b= 0.0001income+19.827credit\ card\ ins-8.314sex+0.415age+17.691$$

با این معادله می توان نتایج Life Insurantee promotion بدست آورد ، که همانطور که در جدول فوق نشان داده شده نتایج مناسبه شده با متغیر وا بسته هم خوانی زیادی دارد. حال اگر نمونه جدیدی به این صورت داشته باشیم:

$$In\ cone=35k \quad credit\ card\ Insurantee=1 \quad sex=0 \quad age=39$$

با محاسبات احتمال بدست آمده برابر 0.999 می باشد. که این فرد یک کاندیدا را برای بیمه عمر (Life (Insurantee promotion می باشد حالت دیگر اینکه نمونه جدید به صورت :

$$Ineome=35k \quad credit\ card\ Insuran=0 \quad sex=1 \quad age=39$$

باشد در این حالت مقدار احتمالی بدست آمده برابر 0.035 است که نشان می دهد یک مرد ۳۹ ساله که در آمد سالیانه او ۳۵۰۰۰ است و بیمه کارت اعتباری ندارد یک نمونه ضعیف برای داشتن بیمه عمر است.

۳- Bayse classsifire

این مقدار یکی از روشهای ساده یادگیری نظارتی است، که در آن فرض می شود که تمام متغیرهای ورودی به یک اندازه مهم هستند و مستقل از هم می باشند و نیز ا گریکی از شرایط هم برقرار نباشد این روش در شرایطی کاربرد دارد این روش بر اساس تئوری بیز بنا شده است. که این تئوری به صورت زیر است:

که در این جا H متغیر وابسته است و E بوسیله مقدار ویژگی های ورودی تعیین می شود.

Bayse classifier برخلاف اکثر روشهای آماری برای حالتی که مقدار داده یک متغیر ورودی نامعلوم است نیز کاربرد دارد. در ادامه با ارائه یک مثال می توان به توضیح این روش پرداخت. در اینجا نیز ابتدا متغیر خروجی را تعیین می کنیم. فرض کنید که یک سری داده داریم، و با استفاده از آن داده ها می خواهیم برای یک نمونه جدید به شکل زیر، جنسیت فرد را تعیین کنیم:

Magazine promotion=Yes watch promotion=Yes
Life Insurance Promotion=No credit card Insurance=No
Sex=?

اگر بخواهیم این نمونه را با فرمول Bayse classifier بنویسیم داریم:
که برابر است با:

و همچنین باید جنسیت زن نیز محاسبه شود یعنی عبارت زیر همانند روش فوق باید را محاسبه کرد.
با محاسبه این احتمال داریم:
و چون $281\% > 593\%$ بنابراین این احتمال اینکه جنسیت فرد در نمونه جدید مرد باشد بیشتر است. پس احتمال اینکه یک نمونه با این مشخصات مذکر باشد و برابر مونث بودن آن است.
مطلب دیگر اینکه زمانی که مقدار یک احتمال صفر باشد چون احتمال ها در هم ضرب می شوند کل احتمال صفر خواهد شد مثلاً وقتی که باشد مقدار احتمال خواهد شد. Bayes classifier برای رفع این مشکل به یک مقدار k به صورت کسر ضرب در یک احتمال p و یک مقدار k به مخرج اضافه می کند. بدین صورت:

می باشد که k مقداری بین صفر و یک دارد که معمولاً مقدار یک می گیرد و نیز p بستگی به تعداد انتخابهای متغیر خروجی دارد مثلاً اگر متغیر خروجی دو حالتی باشد (yes, No)، مقدار p برابر 0.5 می باشد. نیز همان مقدار های یا است. مثلاً اگر مقدار باشد آنگاه برابر خواهد بود.

علاوه بر این روش (Bayse classifier) برای حالت Missing data نیز کاربرد دارد. یعنی اگر مقدار یکی از مشخصه های ورودی در یک نمونه جدید را نداشته باشیم، در این روش این مشخصه را کلاً حذف می شود.

حالت دیگری که می توان این روش را بکار برد وجود مشخصه هایی با داده های عددی در بین مشخصه های ورودی است. مثال زیر که یک نمونه جدید است در نظر بگیرید:

Magazine promotion=Yes , wateh promotion=Yes

Life insurance promotion=No , credite card Insurance=No, Age=45

یعنی همان نمونه قبلی با این تغییر که سن نیز به مشخصه های ورودی اضافه شده. در اینجا با استفاده از این روش ابتدا باید توزیع مشخصه ورودی تعیین کنیم، که معمولاً فرض می شود که مشخصه از

توزیع نرمال پیروی می کند. مثلاً مورد روبرو را می خواهیم حساب کنیم:

که برابر است با:

(عبارت از مثال قبلی که age جز متغیرهای ورودی نبود، بدست آمده)

در اینجا ابتدا باید را حساب کنیم که با بدست آوردن میانگین و واریانس داده های سن داریم:

با قرار دادن در فرمول توزیع نرمال داریم:

که این مقدار برابر با 0.03 است و به همین ترتیب برای حالت نیز حساب می کنیم.

نتیجه لازم را از داده های خروجی می گیریم. با حساب کردن احتمال قانده بیز داریم:

که در اینجا نیز احتمال مرد بودن بیشتر است.

ابزار رگرسیون خطی : (توضیح تابع LINEST)

برای اجرای رگرسیون خطی می توان از نرم افزار Excel استفاده کرد. در قسمت توابع Excle، تابع

LINEST برای اجرای یک رگرسیون خطی ایجاد شده است. در اینجا به چگونگی کار با این تابع و

استفاده از نتایج بدست آمده اجرای آن می پردازیم.

۱- روش کار با نرم افزار: بعد از باز کردن Excel داده هایی را که می خواهیم بوسیله آنها رگرسیون

خطی را اجرا کنیم، وارد می کنیم و نیز محلی را که می خواهیم داده ها خروجی نشان داده شوند تعیین

می کنیم سپس از منوی function. Insert را انتخاب می کنیم. در قسمت select a cotegory مقوله

statistical را انتخاب می کنیم. با انجام این کار تابع های آماری در قسمت پایین همین پنجره نمایش داده می شود. از بین تابع ها تابع LINEST را انتخاب می کنیم و OK می کنیم.

در پنجره LINEST چهار قسمت وجود دارد که باید پر شوند. در قسمت اول باید ستون متغیر وابسته (y) را مشخص می کنیم. مثلاً اگر داده های شما در ستون E از ردیف ۲ تا ۱۲ هستند، در این قسمت می نویسیم E2:E12، در قسمت دوم ستون متغیرهای مستقل ها را به همین ترتیب مشخص می کنیم. قسمت سوم مقدار عدد ثابت رگرسیون را به ما می دهد، اگر این مقدار True انتخاب کنیم، مقدار عدد ثابت را بر می گرداند و اگر False باشد مقدار عدد ثابت صفر است. در قسمت چهارم (state) اگر عبارت Ture تایپ می شود اطلاعاتی را در مورد نتایج رگرسیون بدست می دهد که مفید است، مثلاً ضریب همبستگی، بین مقدار تخمینی و مقدار واقعی متغیر وابسته، مقدار آماره F و... که با مثال بیشتر توضیح داده می شود و اگر این مقدار False باشد Excel این نتایج را ارائه نمی دهد. پس از آنکه هر چهار قسمت پر شد، با نگه داشتن ctrl+shift و زدن inter (یا ok کردن) نتایج اجرای رگرسیون دیده می شود. برای فهم بهتر این ابزار در زیر مثالی آورده می شود.

مثالی را در نظر بگیرید که ۴ متغیر مستقل و یک متغیر وابسته داده ها مربوط به ساختمان اداری می باشد که با داده های ورودی که می گیریم، قصد داریم ارزش یک ساختمان را به واحد پولی دلار تخمین بزنیم. متغیرهای تا و به شرح زیر می باشند:

متغیر	توضیح
Y	قیمت ساختمان
X ₁	میزان فضای ساختمان
X ₂	تعداد اتاقهای ساختمان
X ₃	تعداد ورودی ها
X ₄	میزان عمری که ساختمان داشته به سال

این داده ها بدین شکل در Excel نوشته می شود.

سپس ناحیه ای را که می خواهید داده ها خروجی در آنجا نوشته شود انتخاب کنید
 سپس با آدرس زیر پنجره تابع LINEST را باز کنید.
 با انتخاب LINEST از مقوله statistical، OK کنید،

تا پنجره LINEST باز شود در قسمت known-y عبارت E2:E12 و در قسمت known-x عبارت
 A2:D12 را تایپ کنید، سپس در قسمت const و stats عبارت True را تایپ کنید.

با نگه داشتن ctrl+shift و زدن ok نتایج در جایی که قبلاً انتخاب کردید، نوشته می شود.

231.8145	2709.2	12618.39	25.5609	56587.02
13.72808	549.07	413.9391	5.617636	12661.69
0.996544	1004.233	#N/A	#N/A	#N/A
432.4997	6	#N/A	#N/A	#N/A
1.74E+09	6050904	#N/A	#N/A	#N/A

۲- توضیح نتایج خروجی:

ردیف اول ضرایب متغیرهای مستقل و عدد ثابت b را نشان می دهد. در واقع همان ها و b در معادله هستند که بصورت برعکس از راست به چپ نوشته شده اند. یعنی از سمت راست ترین عدد مقدار b و عدد بعدی ضریب X_1 که floor space می باشد نشان داده شده و آخرین عدد سمت چپ ضریب X_4 یعنی Age می باشد.

ردیف دوم نتایج خروجی خطای استاندارد هر ضریب و عدد ثابت b را نشان می دهد. که میزان انحراف نتایج هر ضریب از مقدار میانگین آن نشان می دهد و مثلاً عدد 13.72808 میزان انحراف ضریب X_4 را از مقدار میانگین آن نشان می دهد.

داده اول ردیف سوم میزان ضریب همبستگی بین مقدار تخمینی متغیر وابسته و مقدار واقعی این متغیر را نشان می دهد که بین 1- و 1 می باشد و هر چه این ضریب به 1-، 1 نزدیکتر باشد نشان می دهد که معادله رگرسیون پیشگویی خوبی برای مقدارهای واقعی متغیر وابسته است و هر چه به صفر نزدیک تر باشد نشان می دهد که روش رگرسیون خطی نامعتبر است که یک معیار برای فهمیدن اینکه رگرسیون خطی مناسب است یا خیر می باشد. داده دوم در همین سطر خطای استاندارد متغیر وابسته را از میانگین آن نمایش می دهد.

داده اول سطر چهارم یکی از داده های خروجی مفید است و معیاری خوب است برای پی بردن به اینکه آیا رگرسیون خطی مناسب است یا خیر. این مقدار آماره F را بدست می دهد. این آماره بعنوان توزیع نمونه گیری دو متغیر تصادفی مستقل که بر درجه آزادیشان تقسیم شده اند، مورد مطالعه قرار می گیرد. برای تفسیر F باید به دو درجه آزادی دسترسی داشته باشیم، این مقادیر اغلب جدول توزیع F را به دو مقدار V_1 و V_2 تفکیک می کند. مقدار V_1 تعداد متغیر های مستقل می باشد که در اینجا برابر 4 می باشد و V_2 حاصل تفریق تعداد کل دادهها از کل متغیر هانیست 5 تا می باشد (تا مستقل و یک وابسته) پس V_2 برابر است با $6 = 11 - 5$ در اینجا از سطح 5% نیز برای آزمون آماره F استفاده می شود. با این توضیح می توان مقدار $F_{0.05, 4, 11}$ را از جدول توزیع F که در کتابهای آمار موجود می باشد بدست آورد. خلاصه مطلب اینکه اگر مقدار اول داده اول سطر چهارم از مقدار $F_{0.05, V_1, V_2}$ که از جدول بدست می آید بزرگتر باشد، معادله رگرسیون به دقت توانسته ارزش متغیر وابسته را تعیین کند یا به عبارتی رابطه مشاهده شده بین متغیر های مستقل و متغیر وابسته اتفافی نیست و اگر این مقدار کوچکتر باشد این روش معتبری برای داده ها نیست، یا به صورت آماری می گوئیم فرض صفر را مبنی بر اتفافی بودن رابطه متغیر های مستقل و متغیر وابسته را نمی توان رد کرد.

داده دوم در همین سطر درجه آزادی V_2 را نشان می دهد. در سطر آخر نیز نشان دهنده مجموع مربعات رگرسیون و مجموع مربعات باقیمانده است که از لحاظ آماری فرمول آن بدین شکل است:

که میانگین است و مقدار هر داده می باشد.