

۱۳۰۰



دانشگاه اصفهان

دانشکده فنی و مهندسی

گروه مهندسی کامپیوتر

پایان نامه ی کارشناسی ارشد رشته ی مهندسی کامپیوتر گرایش نرم افزار

بهبود روش های دسته بندی در داده کاوی با استفاده از دانش گذشته

استاد راهنما:

دکتر احمد برآنی

پژوهشگر:

مهدی مرادیان

۱۳۸۸/۱۰/۲۷

خرداد ماه ۱۳۸۸

اطلاعات مدرک علمی برآید
توسط مدرک

۱۳۰۰۳۷

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

پایان نامه
گرایش کامپیوتر
رشته مهندسی
تصویبات تکمیلی دانشگاه اصفهان



دانشگاه اصفهان

دانشکده فنی و مهندسی

گروه مهندسی کامپیوتر

پایان نامه ی کارشناسی ارشد رشته ی مهندسی کامپیوتر گرایش نرم افزار
آقای مهدی مرادیان تحت عنوان

بهبود روش های دسته بندی در داده کاوی با استفاده از دانش گذشته

در تاریخ ۱۳۸۸/۳/۲۵ توسط هیأت داوران زیر بررسی و با درجه عالی به تصویب نهایی رسید.

امضا

با مرتبه ی علمی استادیار

دکتر احمد برآنی

امضا

با مرتبه ی علمی استادیار

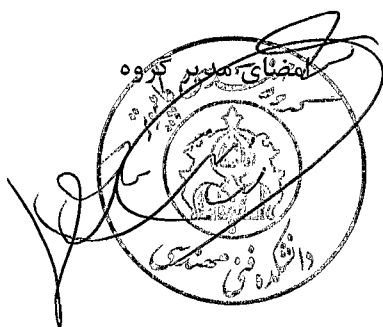
دکتر بهروز ترک لادانی

امضا

با مرتبه ی علمی استادیار

دکتر مازیار پالهنگ

امضا



با تشکر از

استاد ارجمندم جناب آقای دکتر براآنی که
همواره از مساعدت های ایشان در به انجام
رساندن این تحقیق بهره مند بوده ام.

و

فرهیختگان محترم هیئت داوران و اساتید
گرانقدر که داشته هایم از دانش، مرهون
زحمات و راهنمایی های ایشان است.

تقدیم به پیشگاه مقدس حضرت ولی عصر (عج)

تقدیم به پدر ارجمندم

که بزرگوارانه در حفظ و تربیت من کوشید و عزت و
همت والای خود را سرلوحه دفتر زندگیم قرار داد و
اقتدار من حاصل زحمت بیکران اوست.

تقدیم به مادر عزیزم

که خورشید مهربانیش هیچگاه غروب نمی کند و باران
رحمتش همیشگی است. او که دریای عشق و ایثار است
و از دعای اوست که امروز به اینجا رسیده ام.

و

همسر

عزیزم

چکیده

الگوریتم KNN یکی از بهترین و پرکاربردترین الگوریتم های دسته بندی است که از آن استفاده ی گسترده ای در کاربردهای مختلف می شود. یکی از مشکلات این الگوریتم، تأثیر یکسان همه ی خصیصه ها در محاسبه ی فاصله ی رکورد جدید با رکوردهای موجود در پایگاه داده های آموزشی می باشد، در صورتی که برخی از این خصیصه ها برای عمل دسته بندی کم اهمیت ترند. این امر باعث گمراهی روند دسته بندی و کاهش دقت الگوریتم دسته بندی می شود. در این تحقیق به استخراج نوع خاصی از قوانین وابستگی می پردازیم که سمت چپ آن ها فقط یک قلم وجود دارد و سمت راست آن ها برچسب دسته وجود دارد. سپس با بررسی و آنالیز این قوانین وابستگی و ترکیب آن ها با الگوریتم KNN دو الگوریتم دسته بندی جدید پیشنهاد می کنیم. در الگوریتم اول یعنی الگوریتم نزدیک ترین k تایی همسایگی مبتنی بر وابستگی پویا¹ بر اساس مقادیر ویژگی های رکورد جدید به صورت پویا به ویژگی های مختلف وزن اختصاص می دهیم، یعنی هر بار که یک رکورد جدید وارد می شود، وزن ویژگی ها برای محاسبه ی فاصله تغییر می کند. اما در الگوریتم دوم یعنی الگوریتم نزدیک ترین k تایی همسایگی مبتنی بر وابستگی ایستا² بدون توجه به مقادیر ویژگی های رکورد جدید و به صورت ایستا به ویژگی های مختلف وزن اختصاص می دهیم. در الگوریتم دوم وزن همه ی ویژگی ها برای همه ی رکوردهای جدیدی که وارد سیستم می شوند، ثابت است. بعد با توجه به وزن خصیصه های مختلف و با استفاده از فرمول محاسبه ی فاصله ی مانهتن بین رکوردها به دسته بندی بر اساس الگوریتم KNN می پردازیم و با این عمل دقت دسته بندی الگوریتم KNN را افزایش می دهیم. مقایسه ی نتایج ارزیابی این الگوریتم با هفت الگوریتم دیگر دسته بندی بر روی پایگاه داده ی مختلف، بهبود قابل توجه دقت دسته بندی توسط این الگوریتم را نشان می دهد.

کلیدواژه: داده کاوی، دسته بندی، الگوریتم نزدیک ترین k تایی همسایگی، قوانین وابستگی، وزن دهی به خصیصه ها

¹ D_KNNBA (Dynamic-K-Nearest-Neighbor)

² S_KNNBA (Static-K-Nearest-Neighbor)

فهرست مطالب

صفحه

عنوان

فصل اول: کلیات مسئله

- ۱-۱ شرح مسئله..... ۱
- ۲-۱- الگوریتم ارائه شده..... ۳
- ۳-۱- نتایج به دست آمده..... ۴

فصل دوم: داده کاوی

- ۱-۲ مقدمه..... ۶
- ۲-۲ داده کاوی چیست؟..... ۷
- ۳-۲ فرایند کشف دانش (KDD)..... ۸
- ۴-۲ معماری سیستم پایگاه داده ها..... ۱۰
- ۵-۲ روش های کاویدن داده..... ۱۳
- ۱-۵-۲ استخراج الگوهای تکرار شونده و قوانین وابستگی..... ۱۳
- ۲-۵-۲ دسته بندی و پیشگویی..... ۱۴
- ۳-۵-۲ خوشه بندی..... ۱۵
- ۴-۵-۲ آنالیز بخش های جدا..... ۱۶
- ۶-۲ الگوهای جذاب..... ۱۶
- ۷-۲ اولیه های داده کاوی..... ۱۸
- ۸-۲ مسائل متفرقه در داده کاوی..... ۲۰
- ۹-۲ نتیجه گیری..... ۲۱

فصل سوم: دسته بندی

- ۱-۳ مقدمه..... ۲۲
- ۲-۳ دسته بندی چیست؟..... ۲۲
- ۱-۲-۳ کاربردهای دسته بندی..... ۲۳
- ۲-۲-۳ چگونه دسته بندی کنیم؟..... ۲۴

۲۵ ۳-۲-۳ انتخاب داده های آموزشی و داده های آزمایشی
۲۶ ۴-۲-۳ تفاوت دسته بندی با پیشگویی
۲۷ ۳-۳ آماده سازی داده ها برای دسته بندی
۲۸ ۴-۳ الگوریتم های دسته بندی
۲۹ ۱-۴-۳ الگوریتم درخت تصمیم
۲۹ ۱-۱-۴-۳ مثال
۳۱ ۲-۱-۴-۳ الگوریتم پایه ای درخت تصمیم
۳۲ ۳-۱-۴-۳ معیار انتخاب خصیصه ها
۳۳ ۱-۳-۱-۴-۳ معیار بهره دهی اطلاعاتی
۳۴ ۲-۳-۱-۴-۳ معیار نسبت سود
۳۵ ۲-۴-۳ دسته بندی به وسیله ی آنالیز قوانین وابستگی
۳۶ ۱-۲-۴-۳ الگوریتم CBA
۳۶ ۲-۲-۴-۳ الگوریتم CMAR
۳۹ ۳-۲-۴-۳ الگوریتم CPAR
۴۰ ۵-۳ نتیجه گیری
فصل چهارم: الگوریتم KNNBA	
۴۲ ۱-۴ مقدمه
۴۲ ۲-۴ الگوریتم نزدیک ترین K تایی همسایگی
۴۶ ۱-۲-۴ محاسبه ی فاصله ی دور کورد
۴۷ ۲-۲-۴ نرمالسازی مقادیر ویژگی های رکوردها
۴۸ ۳-۲-۴ تابع ترکیب
۵۰ ۴-۲-۴ وزن دهی به خصیصه ها
۵۶ ۵-۲-۴ محاسبه ی مقدار k
۵۷ ۶-۲-۴ پایگاه داده ی مناسب
۵۷ ۷-۲-۴ پیچیدگی زمانی الگوریتم

۵۸ ۸-۲-۴ ترکیب KNN با دیگر الگوریتم های دسته بندی
۵۹ ۳-۴ استخراج قوانین وابستگی
۶۰ ۱-۳-۴ مثال سبد خرید
۶۱ ۲-۳-۴ داده های تراکنشی و جدوال گرا
۶۲ ۳-۳-۴ Apriori الگوریتم
۶۴ ۴-۳-۴ پیچیدگی زمانی الگوریتم Apriori
۶۵ ۴-۴ الگوریتم نزدیک ترین k تایی همسایگی مبتنی بر وابستگی
۶۶ ۱-۴-۴ نحوه کار الگوریتم D_KNNBA
۶۹ ۱-۱-۴-۴ یک مثال از اجرای الگوریتم KNNBA پویا
۷۲ ۲-۱-۴-۴ پیچیدگی زمانی الگوریتم
۷۲ ۱-۲-۱-۴-۴ پیچیدگی زمانی الگوریتم KNN
۷۳ ۲-۲-۱-۴-۴ پیچیدگی زمانی الگوریتم KNNBA
۷۴ ۳-۲-۱-۴-۴ مقایسه پیچیدگی زمانی دو الگوریتم KNN و KNNBA
۷۴ ۲-۴-۴ نحوه کار الگوریتم S_KNNBA
۷۷ ۱-۲-۴-۴ یک مثال از اجرای الگوریتم S_KNNBA
۸۰ ۲-۲-۴-۴ پیچیدگی زمانی الگوریتم S_KNNBA
۸۰ ۵-۴ نتیجه گیری
فصل پنجم: پیاده سازی و ارزیابی الگوریتم های KNNBA	
۸۱ ۱-۵ مقدمه
۸۱ ۲-۵ ارزیابی الگوریتم های KNNBA
۸۲ ۱-۲-۵ مقایسه ی دقت دسته بندی الگوریتم های KNN ، D_KNNBA ، S_KNNBA
۹۲ ۲-۲-۵ مقایسه ی الگوریتم KNNBA با دیگر الگوریتم ها
۹۷ ۳-۵ نتیجه گیری
۹۸ فصل ششم: نتایج و راهکارهای آینده
۱۰۰ منابع و مآخذ

فهرست شکل ها

صفحه	عنوان
۸	شکل ۱-۲ حضور علوم مختلف در فرایند داده کاوی.....
۹	شکل ۲-۲ فرایند کشف دانش از داده ها.....
۱۲	شکل ۳-۲ اجزای مختلف معماری سیستم پایگاه داده ها.....
۳۰	شکل ۱-۳ اطلاعات مشتری های مثال ۳-۴-۱-۱.....
۳۱	شکل ۲-۳ درخت تصمیم مثال ۳-۴-۱-۱.....
۳۲	شکل ۳-۳ الگوریتم پایه ای تولید درخت تصمیم از مجموعه تاپل های یادگیری.....
۴۳	شکل ۱-۴ نمودار بیماران و داروهای تجویزی برای هر یک.....
۴۴	شکل ۲-۴ مشاهده ی سه بیمار جدید در نمودار موجود.....
۴۵	شکل ۳-۴ نمودار مربوط به بیمار ۲ و همسایه هایش.....
۴۵	شکل ۴-۴ نمودار مربوط به بیمار ۳ و همسایه هایش.....
۴۶	شکل ۵-۴ محاسبه ی فاصله ی اقلیدسی بین دو رکورد.....
۶۴	شکل ۶-۴ اجرای الگوریتم Apriori بر روی جدول ۴-۱.....
۶۵	شکل ۷-۴ الگوریتم پایه ای KNN.....
۶۹	شکل ۸-۴ الگوریتم D_KNNBA.....
۷۷	شکل ۹-۴ الگوریتم S_KNNBA.....
۸۳	شکل ۱-۵ دقت الگوریتم KNNBA با تعداد همسایه های متفاوت.....
۸۴	شکل ۲-۵ دقت الگوریتم KNNBA با تعداد همسایه های متفاوت با حذف فیلد کلید از پایگاه داده
۸۵	شکل ۳-۵ دقت الگوریتم KNNBA با حدود آستانه ی ضریب اطمینان متفاوت.....
۸۶	شکل ۴-۵ دقت الگوریتم KNNBA با حدود آستانه ی پشتیبانی متفاوت.....
۸۷	شکل ۵-۵ مقایسه ی دقت دو روش با n های مختلف.....
۹۵	شکل ۶-۵ نمودار دقت دسته بندی الگوریتم ها بر روی ۱۵ پایگاه داده از UCI.....

فهرست جدول ها

صفحه	عنوان
۶۰	جدول ۱-۴ سبد خرید ۴ مشتری.....
۶۹	جدول ۲-۴ داده های مثال های ۱-۴-۴ و ۱-۲-۴-۴.....
۷۰	جدول ۳-۴ اقلام پررخداد استخراج شده از جدول ۲-۴.....
۷۸	جدول ۴-۴ گروه های استخراج شده توسط الگوریتم S_KNNBA از داده های جدول ۲-۴.....
۸۲	جدول ۱-۵ مقایسه ی دقت سه روش با k های مختلف بر روی پایگاه hayes.....
۸۸	جدول ۲-۵ مشخصات پایگاه داده های مورد استفاده از UCI.....
۸۹	جدول ۳-۵ پارامترهای الگوریتم KNNBA برای اجرا بر روی ۱۵ پایگاه داده از UCI.....
۹۰	جدول ۴-۵ دقت اجرای الگوریتم های KNN,D_KNNBA,S_KNNBA بر روی ۱۵ پایگاه داده از UCI.....
۹۴	جدول ۵-۵ دقت دسته بندی الگوریتم ها بر روی ۱۵ پایگاه داده از UCI.....

فصل اول

کلیات مسئله

۱-۱ شرح مسئله

استفاده از داده کاوی به دلیل حجم زیاد داده ها و همچنین میزان کم دانش و اطلاعاتی که بشر با توجه به این حجم داده در اختیار دارد، ضروری به نظر می رسد [۱] - [۲]. داده کاوی فرایندی است که با استفاده از تکنیکهای هوشمند، دانش را از مجموعه ای از داده ها استخراج می کند. در داده کاوی از سه روش اصلی استخراج قوانین وابستگی^۱، خوشه بندی^۲ و دسته بندی^۳ استفاده می شود [۱].

دسته بندی نوعی از آنالیز داده ها می باشد که رکوردها را به وسیله ی الگوهای از پیش استخراج شده^۴، بین دسته های از پیش تعریف شده تقسیم می کند و هر رکورد را با توجه به خصیصه هایش^۵ در یکی از این دسته ها قرار می دهد. دسته ها بر اساس بررسی داده های موجود و با توجه به شرایط و نیاز محیط (تعریف مسئله) تعریف می شوند و به هر دسته معنای خاص و برچسب دسته ی خاص اختصاص داده می شود. برای مثال مشتری های یک بانک را می توان به دو دسته ی مشتری های خوش حساب و مشتری های بد حساب تقسیم کرد که برچسب دسته ی اول، مشتری خوش حساب و معنای آن برای استفاده کننده از سیستم در بانک

^۱ Association rules

^۲ Clustering

^۳ Classification

^۴ در ساده ترین حالت این الگوها تعدادی قوانین اگر... آنگاه... (if... then...) هستند.

^۵ Attribute

مشخص است [۱]. دسته بندی کاربردهای متعددی دارد، به عنوان مثال از دسته بندی برای تشخیص بیماری‌ها^۳، بازاریابی انتخابی^۴، تصویب اعتبار^۵ و غیره استفاده می‌شود. الگوها توسط الگوریتم‌های مختلف دسته بندی از جمله الگوریتم‌های شبکه عصبی، درخت تصمیم، ژنتیک و غیره [۱] - [۵] - [۶] - [۷] - [۸] - [۹] - [۱۰] - [۱۱] استخراج می‌شوند. از الگوهای به دست آمده برای فهم و تشریح داده‌های موجود و برای پیش بینی نحوه رفتار موارد جدید استفاده می‌شود.

البته دسته‌ی دیگری از الگوریتم‌های دسته بندی نیز وجود دارند، که به جای استخراج الگو به منظور دسته بندی تا زمان ورود یک داده‌ی جدید به سیستم کاری انجام نمی‌دهند، هر بار پس از ورود داده‌ی جدید با مقایسه‌ی آن داده و داده‌های موجود در سیستم به دسته بندی داده‌ی جدید می‌پردازند.

الگوریتم نزدیک ترین k تایی همسایگی^۹ یکی از این الگوریتم‌ها است که پس از ورود رکورد جدید به سیستم، k تا از مشابه ترین رکوردهای موجود را به عنوان نزدیک ترین همسایه‌های این رکورد تشخیص داده و بر اساس برچسب دسته‌ی این همسایه‌ها به رکورد جدید یک برچسب دسته اختصاص می‌دهد. در الگوریتم نزدیک ترین k تایی همسایگی، هر بار پس از ورود یک رکورد جدید به سیستم، آن را در یک فضای n بعدی (n تعداد خصیصه‌های رکورد است.) با رکوردهای موجود مقایسه می‌کند و براساس معیار فاصله‌ی اقلیدسی k تا از نزدیک ترین رکوردها به این رکورد جدید را به عنوان همسایه‌هایش انتخاب می‌کند. سپس در ساده ترین حالت برچسب دسته‌ی آن که اکثریت همسایه‌ها متعلق به آن دسته باشند را استخراج می‌کند. استفاده از این الگوریتم در پایگاه داده‌های بزرگ به دلیل حجم زیاد عملیات، کارایی زیادی ندارد ولی در پایگاه داده‌های کوچک یکی از بهترین الگوریتم‌های دسته بندی است.

یکی از مشکلات این الگوریتم این است که وقتی رکورد جدید را با رکوردهای موجود در فضای n بعدی مقایسه می‌کنیم به همه‌ی ویژگی‌ها به یک اندازه توجه می‌کنیم، در صورتی که ویژگی‌های مختلف برای دسته بندی رکوردها اهمیت یکسان ندارند. برخی از ویژگی‌ها هیچ نقشی در دسته بندی ندارند، برخی اهمیت

^۶ با توجه به علائم یک بیماری، الگویی استخراج می‌کند و با آن الگو افراد را به دو دسته، افراد دارای بیماری و افراد فاقد بیماری تقسیم

می‌کند.

^۷ با توجه به ویژگی‌های مشتری‌ها، الگویی استخراج می‌کند و با آن الگو افراد را به دسته‌هایی با سلاقی خرید مختلف تقسیم می‌کند و در هنگام ارسال آگهی‌های تبلیغاتی برای هر دسته مشتری آگهی‌های مناسب را ارسال می‌کند.

^۸ مشتری‌های یک بانک را با توجه به ویژگی‌هایشان از نظر گردش مالی، خوش حسابی و غیره به چند تقسیم می‌کند و به هر دسته تا سقف مشخصی اعتبار اختصاص می‌دهد.

^۹ k-Nearest-Neighbors

کمترو برخی دیگر اهمیت زیادی در دسته بندی آن رکورد جدید دارند. به همین دلیل برخی از تحقیق هایی که سعی در بهبود الگوریتم نزدیک ترین k تایی همسایگی داشته اند، به جای ارزش دومی یکسان به همه ی ویژگی ها، به ویژگی ها، وزن های مختلف نسبت می دهند و به جای استفاده از فاصله ی اقلیدسی از فاصله ی مانهتن استفاده می کنند. این کار باعث کاهش تأثیر منفی ویژگی های غیر مرتبط با روند دسته بندی در محاسبه ی فاصله ی بین رکوردها می شود و در عمل دقت دسته بندی را افزایش می دهد.

یکی دیگر از روش های داده کاوی، استخراج قوانین وابستگی است. در این روش، الگوهای بین داده ها به صورت قوانینی بیان می شود. به عنوان مثال در پایگاه داده های مشتری های بانک ممکن است قانونی به شکل زیر قابل استخراج باشد:

Credit = good AND Age = middle \rightarrow Credit - rating = good

یک قانون وابستگی یک گزاره ی منطقی به شکل $A \rightarrow B$ است که در آن A و B دو مجموعه قلم هستند. که اشتراکشان تهی است. هر قلم ارتباطی بین یک ویژگی و مقدارش را بیان می کند.

A پیش فرض و B نتیجه ی قانون است. در این قانون بیان می شود که مشتری های میان سال که حقوق آنها خوب است، مشتری های خوبی برای بانک هستند. کارایی هر قانون وابستگی توسط دو پارامتر پشتیبانی و ضریب اطمینان سنجیده می شود. پشتیبانی قانون $A \rightarrow B$ ، برابر درصدی از رکوردها است که A و B با هم در آن ها وجود داشته باشند و ضریب اطمینان این قانون، برابر درصدی از رکوردهای شامل A است که شامل B نیز هستند.

در این تحقیق می خواهیم با استفاده از قوانین وابستگی برای هر ویژگی در الگوریتم KNN وزن به دست آوریم و دقت الگوریتم KNN را افزایش دهیم.

۱-۲- الگوریتم ارائه شده

در این تحقیق کلیه ی قوانین وابستگی به صورت $A \rightarrow B$ که در آنها A تنها شامل یک قلم و B بیانگر یک برجسب دسته است، استخراج می شوند. سپس ضریب اطمینان هر قانون و پشتیبانی قلم A را در آن قانون استخراج می کنیم. حال کلیه ی قوانینی که قلم A در آنها مربوط به یک ویژگی یکسان می شود را در یک گروه قرار می دهیم. واضح است که با این کار سمت چپ کلیه ی قوانین هر گروه مثل G_i فقط اقلام مرتبط با یک ویژگی مثل A_i وجود دارد و بنابراین هر گروه نماینده ی یکی از ویژگی ها است.

سپس برای هر گروه دو پارامتر ضریب اطمینان گروه برابر حداکثر مقدار ضرایب اطمینان قوانین مختلف موجود در یک گروه و پشتیبانی برابر حداکثر مقدار پشتیبانی قلم های مختلف موجود در سمت چپ قوانین آن گروه را تعریف می کنیم.

ایده ی اصلی برای تشخیص اهمیت هر ویژگی در دسته بندی این است که چقدر مقادیر آن ویژگی تکرار شده اند و از طرف دیگر تکرار مقادیر آن ویژگی چقدر با تعیین بر حسب کلاس برای رکوردهای آن پایگاه داده ها مرتبط است. بر همین اساس مقادیر دو پارامتر ضریب اطمینان گروه و پشتیبانی گروه را با مقادیر حد آستانه ی تعریف شده در هر پایگاه داده مقایسه می کنیم، ویژگی هایی که مقادیر حداقل یکی از این دو پارامتر در آن ها از حد آستانه ی تعریف شده پایین تر باشد، وزن ϕ خواهند داشت. ($w_i = \phi$) و بقیه ی ویژگی ها براساس پشتیبانی گروه آنها مانند فرمول ۱-۱ وزن می گیرند.

$$w_i = \frac{1}{1 - (G - SUP_i)} \quad \text{فرمول ۱-۱}$$

حال با توجه به این وزن ها و با استفاده از فرمول محاسبه ی فاصله ی ۱-۲ به جای فرمول فاصله ی اقلیدسی، الگوریتم KNN را اجرا می کنیم. این وزن های جدید، باعث بهبود الگوریتم KNN می شوند.

$$\begin{aligned} x_1 &= (x_{11}, x_{12}, \dots, x_{1n}) \\ x_2 &= (x_{21}, x_{22}, \dots, x_{2n}) \end{aligned} \quad \text{فرمول ۲-۱}$$

$$dist(x_1, x_2) = \sqrt{\sum_{i=1}^n w_i * (x_{1i} - x_{2i})^2}$$

۱-۳- نتایج به دست آمده:

الگوریتم ارائه شده از زوایای مختلفی با الگوریتم KNN بر روی پایگاه داده های موجود در UCI [۱۲] مورد ارزیابی قرار گرفت. در فصل ۴ ابتدا به بررسی تاثیر اعداد K ی مختلف بر روی دقت الگوریتم ارائه شده و الگوریتم KNN می پردازیم و مشاهده می کنیم که برای مقادیر مختلف K الگوریتم ارائه شده بهبود دقت را نسبت به الگوریتم KNN نشان می دهد. از طرف دیگر برای مقادیر مختلف n، روش ارزیابی n-Cross-fold-validation را در دو الگوریتم KNN و الگوریتم ارائه شده اجرا کرده و مشاهده می کنیم که مقادیر مختلف n نیز نمی تواند تاثیر چندانی در کیفیت عملکرد الگوریتم ارائه شده نسبت به الگوریتم KNN داشته باشد و در پایان الگوریتم ارائه شده را با ۷ الگوریتم دیگر دسته بندی در نرم افزار Weka [۱۳] بر روی ۱۵

پایگاه داده از UCI اجرا کرده و مشاهده می کنیم که این الگوریتم، متوسط دقت بهتری را نسبت به دیگر الگوریتم ها ارائه می کند.

در فصل ۲ به بررسی اجمالی مفهوم داده کاوی، روش های مختلف آن و غیره می پردازیم. در فصل ۳ به بیان روش های مختلف دسته بندی و مسائل پیرامون آن می پردازیم. در فصل ۴ ابتدا الگوریتم KNN را با ذکر یک مثال بیان کرده و الگوریتم های مختلفی که قصد در بهبود این الگوریتم داشته اند را بیان می کنیم و در ادامه به بیان مفهوم قوانین وابستگی و نحوه استخراج آن ها از داده ها می پردازیم، سپس دو الگوریتم جدید براساس استفاده از نوع خاصی از قوانین وابستگی در الگوریتم KNN پیشنهاد می کنیم. در فصل ۵ به ارزیابی الگوریتم های ارائه شده در این تحقیق از جنبه های مختلف و مقایسه ی آن ها با دیگر الگوریتم های دسته بندی می پردازیم.

فصل دوم

داده کاوی

۲-۱ مقدمه

از سال ۱۹۶۰ که سیستمهای پایگاه داده ها، جای سیستمهای پردازش فایل را گرفتند [۱]، با سرعت این سیستمها رو به پیشرفت بودند و از طرف دیگر کاربرد آنها در عرصه های مختلف صنعت و اقتصاد و سیاست و غیره با سرعت شگفت انگیزی رو به افزایش بود. این رشد روز افزون داده ها و از طرف دیگر، عدم توانایی بشر برای درک و بررسی همه داده ها باعث شد که بیشتر داده ها به آرشیو سپرده شوند و در تصمیم گیریها فقط از داده های کمی (که به نظر مدیر مهم تر بودند) استفاده شود، این درحالی بود که شاید داده های موجود در آرشیو دارای اطلاعات مفید زیادی بودند.

چون نیاز ما در اختراع است. محققان شروع به کار برای به دست آوردن فرایندی کردند که بتواند از میان داده های زیاد، دانش ارزشمند را به دست آورد. آنها به دنبال یک مکانیزم خود کار بودند که بدون دخالت بشر بتواند در میان حجم عظیم داده ها الگوهای مفید، جدید و قابل فهم را استخراج کند. اولین بار دکتر شاپیرو^۱ در سال ۱۹۸۹ یک کارگاه کوچک^۲ برای داده کاوی^۳ ایجاد کرد ولی از آن سال به بعد به مرور زمان آن کارگاه کوچک تبدیل به کنفرانسهای مهم و بزرگ کشف دانش در پایگاه داده ها^۴ شد. داده کاوی تا آنجا گسترش

^۱ Shapiro

^۲ KDD-89: IJCAI-89 workshop on Knowledge Discovery in Databases August 20, 1989, Detroit MI, USA

^۳ Datamining

^۴ KDD

یافت که اگر امروز در موتور جستجوی گوگل، داده کاوی را جستجو کنیم بیش از ۱۰۰۰۰۰۰۰ صفحه را پیدا می کند که نشان دهنده گسترش سریع و زیاد داده کاوی در صنعت، تجارت و غیره است.

۲-۲ داده کاوی چیست؟

داده کاوی استخراج دانش ارزشمند از بین حجم زیادی داده است [۱]. اما به نظر می رسد که این نام چندان مناسب نیست. وقتی طلا را از سنگ آن استخراج می کنیم آنگاه می گوئیم به کاویدن طلا پرداختیم نه کاویدن سنگ، بنابراین باید از دانش کاوی^۱ به جای داده کاوی استفاده کنیم. محققان مختلف برای فرایند داده کاوی نام های مختلفی را استفاده کرده اند، از جمله عصاره گیری دانش^۲، آنالیز داده / الگو^۳، باستان شناسی داده^۴ و غیره. داده کاوی را می توان به درختی تشبیه کرد که در علوم مختلف ریشه دارد. از جمله این علوم می توان به تکنولوژی پایگاه داده ها، یادگیری ماشین، آمار، شخصی سازی^۵ و بصری سازی^۶ اشاره کرد (شکل ۲-۱). در ادامه جای پای این علوم را در فرایند داده کاوی مشاهده خواهیم کرد.

برخی از محققان کل فرایند کشف دانش از داده ها^۷ را داده کاوی می گویند ولی برخی دیگر داده کاوی را فقط یکی از مراحل فرایند کشف دانش می دانند.

برای طراحی یک دوره حیات^۸ برای فرایند داده کاوی شرکت های مختلف پیشنهاد های متفاوتی داده اند. معروف

ترین آنها CRISP-DM [۲] از شرکت SPSS و SEMMA از شرکت SAS است ولی در این قسمت به

بررسی فرایندی می پردازیم که در بیشتر کتاب های آکادمیک از این فرایند استفاده می شود.

^۱ Knowledge mining

^۲ Knowledge extraction

^۳ data/pattern Analysis

^۴ Data Archaeology

^۵ privacy

^۶ visualization

^۷ Knowledge Discovery from Data(KDD)

^۸ life cycle