



دانشکده فنی

پایان نامه کارشناسی ارشد

رشته مهندسی صنایع - صنایع

عنوان پایان نامه :

دسته بندی داده ها به روش ماشین های بردار پشتیبان با توابع هدف چندگانه

نگارش : علی ندائی

استاد راهنما : دکتر علی محمد احمدوند

تابستان ۹۱

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

کلیه حقوق این پژوهش متعلق به دانشگاه شاهد بوده و هر گونه استفاده از مطالب و محتویات آن منوط به ذکر منبع و هر گونه کپی برداری از تمام یا بخشی از اثر تنها با کسب مجوز کتبی از دانشگاه مذکور امکان پذیر می باشد.

صورت جلسه دفاع

پدر و مادر عزیزم

تشکر و قدردانی:

از خدای بزرگ برای سلامتی و پشتکاری که به من عنایت فرمود تا توانایی پشت سر گذاشتن فراز و فرودها و به سرانجام رساندن بخشی از اهداف زندگیم را بدست آورم سپاسگزارم.

از پدر و مادر عزیزم و سایر اعضای خانواده ام کمال تشکر را دارم و از محضر خدای متعال برایشان آرزوی توفیق روز افزون می نمایم. این اثر ناچیز را به پاس سال ها تلاش پدر و مادرم به محضرشان تقدیم می کنم.

همچنین از جناب آقای دکتر احمدوند سپاسگذارم که با راهنمایی ها و کمک های بی دریغشان، نه تنها در این پایان نامه و مقطع تحصیلی و نه فقط از جهت علمی، بلکه برای تمام طول عمرم درس های فراوانی به من آموختند. از خدای متعال برای ایشان و خانواده شان موفقیت و شادکامی روز افزون خواستارم.

چکیده:

ماشین بردار پشتیبان روشی بر پایه مفاهیم تئوری یادگیری آماری می باشد که در آن با استفاده از یک مدل برنامه ریزی ریاضی ابرصفحه ای برای دسته بندی داده ها یافته می شود. در این پژوهش مدلی برای ماشین بردار پشتیبان با نرخ خطای چندگانه ارائه می گردد که در آن خطای دسته بندی در دسته اول و دوم با هم برابر نمی باشند. مدل مذکور در حالت خاص معادل ماشین بردار پشتیبان معمولی بوده و هرگونه تفاوت نرخ خطا، بین دو مدل ایجاد اختلاف خواهد نمود.

پس از معرفی و مدلسازی، روش مورد بررسی با استفاده از مجموعه داده هایی تحت آزمایش قرار گرفته و نتایج نهایی از لحاظ دقت دسته بندی و هزینه های ناشی از خطای دسته بندی گزارش خواهند گردید.

واژگان کلیدی - ماشین بردار پشتیبان - دسته بندی - نرخ خطای چندگانه - توابع کرنل - ماشین بردار پشتیبان با حساسیت هزینه

فهرست مطالب

فصل اول: کلیات تحقیق..... ۱۲

۱۳	مقدمه
۱۴	۱ + تعریف مسئله و بیان موضوعات تحقیق
۱۵	۱ ۴ مفروضات تحقیق
۱۶	۱ ۳ ضرورت انجام تحقیق
۱۷	۱ ۴ کاربردهای تحقیق
۱۷	۱ ۵ روش های جمع آوری اطلاعات
۱۸	۱ ۶ نحوه شبیه سازی
۱۸	۱ ۷ ساختار پایان نامه
۱۹	۱ ۸ نوآوری های پایان نامه

فصل دوم: مروری بر ادبیات موضوع..... ۲۰

۲۱	۲-۱ مقدمه
۲۲	۲-۲ مروری بر مفاهیم پایه
۲۲	۲-۲-۱ ابرصفحه و نیم فضا
۲۲	۲-۲-۲ نُرم بردار
۲۳	۲-۲-۳ بردار نرمال ابرصفحه
۲۳	۲-۲-۴ فاصله نقطه از ابرصفحه
۲۴	۲-۲-۵ ابرصفحه متعارفی

۲۴	۳-۲ ماشین بردار پشتیبان
۲۴	۱-۳-۲ مدل داده های تفکیک پذیر خطی
۲۸	۲-۳-۲ مدل داده های تفکیک ناپذیر خطی
۳۰	۳-۳-۲ ماشین بردار پشتیبان با حساسیت هزینه
۳۳	۴-۲ توابع کرنل
۳۶	۵-۲ ماشین بردار پشتیبان چند دسته ای
۳۶	۱-۵-۲ رویکرد یک دسته در مقابل همگی
۳۷	۲-۵-۲ رویکرد یک دسته در مقابل دیگری
۳۷	۳-۵-۲ رویکرد گراف جهت دار بدون دور
۳۸	۴-۵-۲ رویکرد نیم در مقابل نیم

۴۲..... فصل سوم: مدل های پیشنهادی

۴۳	۱-۳ مقدمه
۴۳	۲-۳ تعریف مسئله
۴۴	۳-۳ نرخ خطای چندگانه
۴۵	۴-۳ مدل سازی

۴۹..... فصل چهارم: شبیه سازی و مثال عددی

۵۰	۱-۴ مقدمه
۵۱	۲-۴ الگوریتم نقطه درونی

۵۲	۳-۴ روش حل متوالی برنامه ریزی درجه دو
۵۳	۴-۴ الگوریتم ابتکاری بهبود جواب
۵۴	۵-۴ تحلیل پیچیدگی الگوریتم
۵۴	۶-۴ تعیین ماتریس نرخ خطا
۵۵	۷-۴ حل مدل و نتایج محاسباتی
۶۲	۸-۴ آنالیز حساسیت

۶۸..... فصل پنجم: جمع بندی و نتیجه گیری

۶۹	۱-۵ مقدمه
۶۹	۲-۵ جمع بندی و نتیجه گیری
۷۱	۳-۵ پیشنهادات برای مطالعات آتی

فهرست اشکال

- شکل ۱-۲ تفکیک داده ها در حالت تفکیک پذیر خطی ۲۵
- شکل ۲-۲ تفکیک داده ها در حالت تفکیک ناپذیر خطی ۲۸
- شکل ۳-۲ ماشین بردار پشتیبان غیر خطی ۳۴
- شکل ۴-۲ رویکرد نیم در مقابل نیم ۳۸
- شکل ۱-۴ تکرارهای الگوریتم کارمارکار ۵۱
- شکل ۲-۴ دسته بندی به روش ماشین بردار پشتیبان: (۱) با نرخ خطای چندگانه (۲) با حساسیت هزینه ۶۱
- شکل ۳-۴ مقایسه عملکرد ماشین بردار پشتیبان با حساسیت هزینه و ماشین بردار پشتیبان با نرخ خطای چندگانه ۶۶

فهرست جداول

- جدول ۱-۲ برخی از مدل های توسعه یافته در حوزه ماشین بردار پشتیبان (به ترتیب سال تحقیق) ۳۳
- جدول ۲-۲ برخی از مطالعات صورت گرفته در حوزه مرتبط با توابع کرنل در ماشین بردار پشتیبان (به ترتیب سال تحقیق) ۳۵
- جدول ۳-۲ تحقیقات انجام شده برای ماشین بردار پشتیبان چند دسته ای (به ترتیب سال تحقیق) ۳۹
- جدول ۴-۲ مقایسه روش های دسته بندی داده ها در حالت چند دسته ای ۴۰
- جدول ۱-۴ مجموعه داده های مورد استفاده در حل مدل پیشنهادی ۵۹
- جدول ۲-۴ مقایسه دقت و هزینه خطای مدل های ماشین بردار پشتیبان با ماتریس هزینه ۴-۴ ۶۲
- جدول ۳-۴ مقایسه خطای دسته بندی مدل های ماشین بردار پشتیبان با ماتریس هزینه ۴-۴ ۶۲
- جدول ۴-۴ مقایسه دقت و هزینه خطای مدل های ماشین بردار پشتیبان با $\gamma_{12} = 5$ و $\gamma_{21} = 1$ ۶۳
- جدول ۵-۴ مقایسه خطای دسته بندی مدل های ماشین بردار پشتیبان با $\gamma_{12} = 5$ و $\gamma_{21} = 1$ ۶۳
- جدول ۶-۴ مقایسه دقت و هزینه خطای مدل های ماشین بردار پشتیبان با $\gamma_{12} = 8$ و $\gamma_{21} = 1$ ۶۴
- جدول ۷-۴ مقایسه خطای دسته بندی مدل های ماشین بردار پشتیبان با $\gamma_{12} = 8$ و $\gamma_{21} = 1$ ۶۴
- جدول ۸-۴ مقایسه دقت و هزینه خطای مدل های ماشین بردار پشتیبان با $\gamma_{12} = 10$ و $\gamma_{21} = 1$ ۶۵
- جدول ۹-۴ مقایسه خطای دسته بندی مدل های ماشین بردار پشتیبان با $\gamma_{12} = 10$ و $\gamma_{21} = 1$ ۶۵

فصل اول

کلیات تحقیق

۱-۱ مقدمه

یادگیری ماشینی^۱ یکی از حوزه های میان رشته ای است که هدف اصلی آن یافتن روش های حل و ابداع الگوریتم هایی است که در حل مسائل طبیعی و روزمره راه گشا می باشند. ماشین بردار پشتیبان^۲ از زیر شاخه های یادگیری ماشینی بوده و بر اساس مفاهیم تئوری یادگیری آماری^۳ پایه گذاری شده است. بیشترین محبوبیت ماشین بردار پشتیبان در تشخیص دستخط^۴ و تشخیص صدا^۵ می باشد که از حیث دقت در این حوزه از قویترین شبکه های عصبی نیز پیشی می گیرد. مدل اولیه ماشین بردار پشتیبان توسط وپنیک^۶ (۱۹۹۸) مطرح شد. مدل مذکور یک مدل برنامه ریزی درجه دو^۷ می باشد که در شرایط دو دسته ای به کار می رود. هر چند که این روش در سال های اخیر مورد توجه بسیاری از پژوهشگران بوده است، اما علیرغم این حقیقت، همچنان خلایه های در حوزه مذکور به چشم می خورد.

¹ Machine learning

² Support Vector Machine (SVM)

³ Statistical learning theory

⁴ Handwriting recognition

⁵ Voice recognition

⁶ Vapnik

⁷ Quadratic programming

لذا مدل ماشین بردار پشتیبان در حالت های توسعه یافته تر معمولا با توابع هدف چندگانه مدل می شود. این اصطلاح (ماشین بردار پشتیبان با توابع هدف چندگانه) در ادبیات موضوع مرسوم بوده و به حالت هایی اطلاق می شود که مدل مسئله علاوه بر هدف اصلی، اهداف دیگری را نیز دنبال می کند و انگیزه استفاده از آن نزدیک تر نمودن مدل به دنیای واقعی می باشد. در این فصل در مورد کلیات این تحقیق مطالبی ذکر خواهد شد. همچنین در مورد جنبه های نوآوری طرح مباحثی مطرح می گردید.

۱-۲ تعریف مسئله و بیان موضوعات تحقیق

دسته بندی، به دنبال یافتن مدلی پیشگویانه^۱ می باشد که با استفاده از آن مدل، بتوان در مورد مسائل سازمان در سطوح مختلف تصمیم گیری نمود. این ابزار با توجه به کاربرد قابل ذکر، همواره به عنوان یکی از اصلی ترین تکنیک های موثر در تصمیم گیری در نظر گرفته می شود. روش های زیادی برای دسته بندی معرفی شده اند که مهمترین و پرکاربردترین آن ها در یک مقاله علمی- مروری توسط کاتسیانیتیس^۲ (۲۰۰۷) مورد بررسی قرار گرفته است. به طور خلاصه می توان به درخت تصمیم، شبکه عصبی، الگوریتم های بیزی و ... اشاره نمود.

یکی از نوین ترین و دقیق ترین روش های دسته بندی، ماشین بردار پشتیبان می باشد. در واقع این روش در ابتدا برای دسته بندی داده ها معرفی گردید و مدتی بعد از آن برای کشف الگوهای موجود در میان داده ها برای حالت های توسعه یافته تر مورد استفاده قرار گرفت (وپنیک، ۱۹۹۸). هدف اصلی در روش ماشین های بردار پشتیبان، حداقل نمودن خطای دسته بندی می باشد که این کار با حداکثر نمودن حاشیه دسته بندی صورت می پذیرد. بنابراین آنچه که در مدل ماشین بردار پشتیبان به عنوان تابع هدف ملاحظه می شود، حداکثر نمودن حاشیه است نه حداکثر نمودن دقت دسته بندی.

هدف از انجام این تحقیق یافتن پاسخ سوالات زیر است:

۱- تابع هدف مرسوم در روش ماشین های بردار پشتیبان، نرخ خطای دسته بندی را یکسان در نظر می گیرد. این مطلب ممکن است برخی از واقعیت های موجود را خدشه دار کند. به عنوان مثال یک بانک را در نظر بگیرید که در

¹ Predictive Model

² Kotsiantis

آن می توان مشتریان متقاضی وام را به سه دسته پرخطر، کم خطر و بی خطر تقسیم بندی نمود. اگر یک شخص بی خطر را به اشتباه در دسته پرخطر قرار داده و به شخص مذکور وام تعلق نگیرد ممکن است هزینه تحمیل شده کمتر از حالتی باشد که به یک شخص پرخطر به اشتباه وام تعلق گیرد. بنابراین آیا می توان مدلی با در نظر گرفتن نرخ خطای چندگانه ارائه نمود؟

۲- رویکردهای متفاوتی برای در نظر گرفتن نرخ خطای چندگانه در الگوریتم های دسته بندی معرفی شده اند. از آن جمله می توان به مقاله ارائه شده توسط لی^۱ و همکاران (۲۰۰۹) اشاره نمود. آیا امکان بهبود چنین رویکردهایی برای ماشین بردار پشتیبان وجود دارد؟

۳- در مدل اصلی ماشین بردار پشتیبان، متغیرهای مصنوعی موجود در قیود و تابع هدف به گونه ای تعریف شده اند که اگر مقدار آنها بزرگتر از ۱ باشد، نقطه متناظر دارای خطا در دسته بندی می باشد. رویکردهای مورد بحث تاکنون، نقاطی که متغیر مصنوعی متناظر آنها دارای مقداری بین صفر و یک می باشد را نیز در نظر می گیرند. در حالیکه این نقاط دچار خطای دسته بندی نمی باشند و بنابراین نباید در تعیین محل ابرصفحه دسته بند موثر واقع شوند.

۳-۱ مفروضات تحقیق

- ۱- فرآیند پیش پردازش^۲ بر روی داده ها انجام شده است. بنابراین داده ی مخدوش، ناقص، ناسازگار و... وجود ندارد.
- ۲- مشخصه های^۳ مورد استفاده برای داده ها در دسته بندی موثر می باشند. بنابراین نیاز به حذف هیچکدام از مشخصه ها (بدلیل عدم اهمیت آن ها) نمی باشد.
- ۳- مقادیر موجود در کلیه مشخصه ها برای همه ی داده ها به صورت قطعی اند. بنابراین حالت احتمالی^۴، تصادفی^۵، فازی^۶ و هر حالتی که فرض قطعیت را مختل نماید در میان داده ها وجود ندارد.
- ۴- خطای اندازه گیری و خطای انسانی در میان داده ها وجود ندارد.

¹ Li

² Pre-processing

³ Features

⁴ Probabilistic

⁵ Stochastic

⁶ Fuzzy

۵ - مشخصه ها از هم مستقل اند.

۶ - هزینه ی جمع آوری داده ها (پر نمودن مشخصه ها^۱) در نظر گرفته نمی شود (برابر صفر می باشد).

۷ - داده های مورد استفاده از سایت های مرتبط جمع آوری شده اند. این داده ها به صورت متعادل^۲ می باشند. یعنی تقریباً نیمی از آن ها متعلق به دسته ۱ و نیمی دیگر متعلق به دسته ۲ می باشند. بنابراین نیازی به متعادل سازی داده ها نخواهد بود.

۴-۱ ضرورت انجام تحقیق

به طور کلی از روش های پژوهش عملیاتی به منظور افزایش دقت و ایجاد یک روش قانون محور^۳ در دسته بندی استفاده می شود. از برنامه ریزی ریاضی می توان برای ایجاد مدل هایی با محدودیت های خاص تر استفاده نمود. به عنوان مثال کاریزوسا^۴ و باراگان^۵ (۲۰۰۶) و مالدونادو^۶ و همکاران (۲۰۱۱) با استفاده از تکنیک برنامه ریزی ریاضی محدودیت هایی را به مدل های دسته بندی اضافه نموده اند. هدف از افزودن چنین محدودیت هایی بالا بردن دقت و یا جنبه واقع گرایانه مدل می باشد. به طور کلی دو نوع تقابل بین روش های دسته بندی و روش برنامه ریزی ریاضی وجود دارد:

۱ - روش هایی که بر اساس مفاهیم بهینه سازی و پژوهش عملیاتی توسعه داده شده اند و بدون چنین مفاهیمی بی معنی هستند. مثل ماشین های بردار پشتیبان. در چنین حالت هایی مفاهیم پژوهش عملیاتی از روش مورد استفاده جدا پذیر نیستند.

۲ - روش هایی که مفاهیم بهینه سازی و پژوهش عملیاتی در آن ها به عنوان یک روش جدید به کار رفته است. مثلاً درخت تصمیم که از روش های پژوهش عملیاتی نیز ساخته می شود.

¹ Features

² Balance

³ Rule-based

⁴ Carrizosa

⁵ Baragan

⁶ Maldonado

با توجه به فعالیت های انجام شده در حوزه ماشین بردار پشتیبان به نظر می رسد با ارائه مدل های برنامه ریزی ریاضی ذکر شده در قسمت ۱-۲ می توان قابلیت های مدل ماشین بردار پشتیبان را به عنوان یکی از پرکاربردترین روش های دسته بندی به واقعیت نزدیک تر نمود. روش های توضیح داده شده می توانند کارایی دسته بندی و صحت آن را بهبود بخشند.

۱ ۴ کاربردهای تحقیق

در این تحقیق مفاهیم اولیه در حوزه ماشین بردار پشتیبان به طور کامل توضیح داده شده است. بنابراین این تحقیق می تواند مرجع مناسبی برای مطالعه در حوزه ماشین بردار پشتیبان تلقی گردد. از دیدگاه کاربردی نیز این تحقیق می تواند در کلیه مراکز و سازمان هایی که قصد طبقه بندی داده ها و استفاده از آن ها به منظور پیش بینی را دارند نیز مورد استفاده قرار بگیرد.

۱ ۵ روش های جمع آوری اطلاعات

روش های جمع آوری اطلاعات شامل جستجو در سایت های اینترنتی معتبر علمی، مطالعه مقالات و کتاب ها و مشاوره با استاد راهنما و استفاده از تجربیات و پژوهش های ایشان است. همچنین مشاوره با اساتید این حوزه با نظارت و مشورت استاد راهنما کمک زیادی در گردآوری اطلاعات نموده است. داده های مربوط به شبیه سازی نیز از پایگاه اطلاعاتی مربوط به ماشین بردار پشتیبان به آدرس <http://archive.ics.uci.edu/ml/> استخراج گردیده است. داده های مذکور با فرمت *.txt* دریافت شده و سپس به فرم ماتریسی به نرم افزار *MATLAB* منتقل شده و با فرمت *.mat* مورد تحلیل قرار گرفته اند. داده های مذکور در مقالات متعدد و معتبر به کار گرفته شده است و دسترسی به مقالاتی که از هر یک از داده ها استفاده نموده اند امکان پذیر می باشد. به علاوه به منظور تسهیل بررسی نتایج محاسبتی این تحقیق، نام هر یک از مجموعه داده های^۱ استفاده شده نیز ذکر شده است.

^۱ Datasets

۱ ۶ نحوه شبیه سازی

شبیه سازی و نتایج محاسباتی موجود در این تحقیق با استفاده از نرم افزار *MATLAB R2010a* و الگوریتم های موجود در جعبه ابزار آن انجام شده است. این الگوریتم ها عبارتند از الگوریتم نقطه درونی و روش بهینه سازی متوالی درجه دو که با استفاده از آن مسئله برنامه ریزی درجه دو با محدودیت تبدیل به یک مسئله بدون محدودیت می شود. همچنین پس از حل مسئله با الگوریتم های مطرحه، در فاز دوم با استفاده از یک الگوریتم ابتکاری به بهبود جواب پرداخته شده است. در برخی موارد برای اطمینان از صحت عملکرد نرم افزار، خروجی به دست آمده با خروجی حاصل از سایر منابع و نرم افزارها مقایسه شده است.

۱ ۷ ساختار پایان نامه

در فصل دوم پایان نامه، در ابتدا به تشریح مفاهیم اولیه مورد نیاز مدل ماشین بردار پشتیبان پرداخته شده است. هدف از این قسمت آشنایی خواننده با مدل و نحوه عملکرد ماشین بردار پشتیبان می باشد. مرور ادبیات و پژوهش های مرتبط با این تحقیق در ادامه فصل ۲ معرفی گردیده و مدل هایی که تاکنون ارائه شده اند مورد بحث و بررسی قرار خواهند گرفت.

در فصل سوم نرخ خطای چندگانه^۱ به عنوان یک رویکرد جدید در ماشین بردار پشتیبان مورد بررسی قرار خواهد گرفت. به عبارتی برای این حالت برخلاف حالت های مرسوم، نرخ خطا به صورت یک ماتریس مربعی خواهد بود، بطوریکه مرتبه این ماتریس برابر تعداد دسته های موجود و عناصر قطر اصلی آن همگی برابر صفر می باشند. برای حل این مسئله در فاز اول از الگوریتم های حل مسئله برنامه ریزی درجه دوم (مثل شرط بهینگی کاروش کوهن- تاکر^۲ یا روش لاگرانژین^۳ و یا روش های نقطه درونی^۴ و الگوریتم برنامه ریزی درجه دو متوالی^۵) استفاده شده و در فاز دوم نیز از یک الگوریتم ابتکاری^۶ جواب به دست آمده بهبود داده می شود.

¹ Multiple error rate

² Karush Kohn-Tucker Condition

³ Lagrange

⁴ Interior point

⁵ Sequential quadratic programming

⁶ Heuristic

در فصل چهارم نیز به حل مدل های مذکور و مقایسه آن ها با روش های مرسوم پرداخته شده است. همچنین در مورد روش های حل هر یک از مدل ها توضیحاتی ارائه می شود. مرتبه^۱ زمانی الگوریتم ها نیز مورد تفسیر قرار خواهند گرفت. در انتهای فصل چهارم با تغییر در نرخ خطای دسته بندی و انجام آنالیز حساسیت^۲، مقدار تغییرات در دقت دسته بندی و خطای ایجاد شده در هر دسته مورد بررسی شده و مباحث اصلی خاتمه خواهند پذیرفت. نهایتاً فصل پنجم از این تحقیق به جمع بندی مباحث و پیشنهاداتی برای مطالعات آتی خواهد پرداخت.

۱ A نوآوری های پایان نامه

کلیه مدل های ارائه شده در فصل سوم به همراه نتایج محاسباتی مربوطه در فصل چهارم از جنبه های نوآوری این تحقیق محسوب می شود. تا کنون مدلی برای نرخ خطای چندگانه با در نظر گرفتن نقاطی که دارای خطای دسته بندی می باشند ارائه نشده است. بنابراین مدل ارائه شده در فصل سوم و الگوریتم ابتکاری معرفی شده برای بهبود جواب بدست آمده از جنبه های نوآوری طرح می باشد.

¹ Order

² Sensitivity analysis