





دانشگاه تربیت مدرس
دانشکده علوم پزشکی

پایان نامه

دوره کارشناسی ارشد رشته آمار زیستی

عنوان

شبکه عصبی مصنوعی با مدل خبره-آمیخته و کاربرد آن در
تحلیل داده های پزشکی

دانشجو

مرضیه ابراهیمی

استاد راهنما:

دکتر انوشیروان کاظم نژاد

استاد مشاور:

دکتر محمد غلامی فشارکی

زمستان ۱۳۹۳



تاییدیه اعضای هیات داوران حاضر در جلسه دفاع از
پایان نامه کارشناسی ارشد

خانم مرضیه ابراهیمی رشته آمار زیستی پایان نامه کارشناسی ارشد خود را با عنوان « شبکه عصبی مصنوعی با مدل خبره - آمیخته و کاربرد آن در تحلیل داده های پزشکی » در تاریخ ۱۳۹۳/۱۲/۲۰ ارائه کردند.

بدینوسیله اعضای هیات داوران نسخه نهایی این پایان نامه را از نظر فرم و محتوا تایید کرده و پذیرش آنرا برای تکمیل درجه کارشناسی ارشد پیشنهاد می کنند.

| امضاء | نام و نام خانوادگی | اعضای هیات داوران |
|---|--------------------------|--|
|  | دکتر انوشیروان کاظم نژاد | استاد راهنما |
|  | دکتر محمد غلامی فشارکی | استاد مشاور |
|  | دکتر فرید زایری | استاد ناظر |
|  | دکتر علی اکبر راسخی | استاد ناظر و نماینده تحصیلات تکمیلی |

آیین نامه حق مالکیت مادی و معنوی در مورد نتایج پژوهشهای علمی

دانشگاه تربیت مدرس

مقدمه: با عنایت به سیاست‌های پژوهشی و فناوری دانشگاه در راستای تحقق عدالت و کرامت انسانها که لازمه شکوفایی علمی و فنی است و رعایت حقوق مادی و معنوی دانشگاه و پژوهشگران، لازم است اعضای هیأت علمی، دانشجویان، دانش‌آموختگان و دیگر همکاران طرح، در مورد نتایج پژوهشهای علمی که تحت عناوین پایان‌نامه، رساله و طرحهای تحقیقاتی با هماهنگی دانشگاه انجام شده است، موارد زیر را رعایت نمایند:

ماده ۱- حق نشر و تکثیر پایان‌نامه/ رساله و درآمدهای حاصل از آنها متعلق به دانشگاه می باشد ولی حقوق معنوی پدید آورندگان محفوظ خواهد بود.

ماده ۲- انتشار مقاله یا مقالات مستخرج از پایان‌نامه/ رساله به صورت چاپ در نشریات علمی و یا ارائه در مجامع علمی باید به نام دانشگاه بوده و با تایید استاد راهنمای اصلی، یکی از اساتید راهنما، مشاور و یا دانشجوی مسئول مکاتبات مقاله باشد. ولی مسئولیت علمی مقاله مستخرج از پایان‌نامه و رساله به عهده اساتید راهنما و دانشجو می باشد.

تبصره: در مقالاتی که پس از دانش‌آموختگی بصورت ترکیبی از اطلاعات جدید و نتایج حاصل از پایان‌نامه/ رساله نیز منتشر می‌شود نیز باید نام دانشگاه درج شود.

ماده ۳- انتشار کتاب و یا نرم افزار و یا آثار ویژه (اثری هنری مانند فیلم، عکس، نقاشی و نمایشنامه) حاصل از نتایج پایان‌نامه/ رساله و تمامی طرحهای تحقیقاتی کلیه واحدهای دانشگاه اعم از دانشکده ها، مراکز تحقیقاتی، پژوهشکده ها، پارک علم و فناوری و دیگر واحدها باید با مجوز کتبی صادره از معاونت پژوهشی دانشگاه و براساس آئین‌نامه‌های مصوب انجام شود.

ماده ۴- ثبت اختراع و تدوین دانش فنی و یا ارائه یافته ها در جشنواره‌های ملی، منطقه‌ای و بین‌المللی که حاصل نتایج مستخرج از پایان‌نامه/ رساله و تمامی طرحهای تحقیقاتی دانشگاه باید با هماهنگی استاد راهنما یا مجری طرح از طریق معاونت پژوهشی دانشگاه انجام گیرد.

ماده ۵- این آیین‌نامه در ۵ ماده و یک تبصره در تاریخ ۸۷/۴/۱ شورای پژوهشی و در تاریخ ۸۷/۴/۲۳ در هیأت رئیسه دانشگاه به تایید رسید و در جلسه مورخ ۸۷/۷/۱۵ شورای دانشگاه به تصویب رسیده و از تاریخ تصویب در شورای دانشگاه لازم‌الاجرا است.

«اینجانب مرضیه ابراهیمی دانشجوی رشته آمار زیستی ورودی سال تحصیلی ۱۳۹۱ مقطع کارشناسی ارشد دانشکده علوم پزشکی متعهد می شوم کلیه نکات مندرج در آیین نامه حق مالکیت مادی و معنوی در مورد نتایج پژوهش‌های علمی دانشگاه تربیت مدرس را در انتشار یافته‌های علمی مستخرج از پایان‌نامه / رساله تحصیلی خود رعایت نمایم. در صورت تخلف از مفاد آیین‌نامه فوق‌الاشعار به دانشگاه وکالت و نمایندگی می‌دهم که از طرف اینجانب نسبت به لغو امتیاز اختراع بنام بنده و یا هرگونه امتیاز دیگر و تغییر آن به نام دانشگاه اقدام نماید. ضمناً نسبت به جبران فوری ضرر و زیان حاصله براساس برآورد دانشگاه اقدام خواهم نمود و بدینوسیله حق هرگونه اعتراض را از خود سلب نمودم.»

امضا
تاریخ
۹۴۱۳۱۵

آئین نامه پایان نامه (رساله) های دانشجویان دانشگاه تربیت مدرس

نظر به اینکه چاپ و انتشار پایان نامه (رساله) های تحصیلی دانشجویان دانشگاه تربیت مدرس، مبین بخشی از فعالیت های علمی پژوهشی دانشگاه است. بنابراین به منظور آگاهی و رعایت حقوق دانشگاه، دانش آموختگان این دانشگاه نسبت به رعایت موارد ذیل متعهد می شوند:

ماده ۱ : در صورت اقدام به چاپ پایان نامه (رساله) ی خود، مراتب را قبلاً به طور کتبی به دفتر "دفتر نشر آثار علمی" دانشگاه اطلاع دهد.

ماده ۲ : در صفحه سوم کتاب (پس از برگ شناسنامه)، عبارت ذیل را چاپ کند:

" کتاب حاضر، حاصل پایان نامه کارشناسی ارشد نگارنده در رشته آمار زیستی است که در سال ۱۳۹۳ در دانشکده علوم پزشکی دانشگاه تربیت مدرس به راهنمایی دکتر انوشیروان کاظم نژاد ، مشاوره دکتر محمد غلامی از آن دفاع شده است.

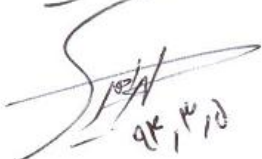
ماده ۳ : به منظور جبران بخشی از هزینه های انتشارات دانشگاه، تعداد یک درصد شمارگان کتاب (در هر نوبت چاپ) را به "دفتر نشر آثار علمی" دانشگاه اهداء کند. دانشگاه می تواند مازاد نیاز خود را به نفع مرکز نشر در معرض فروش قرار دهد.

ماده ۴ : در صورت عدم رعایت ماده ۳، ۵۰٪ بهای شمارگان چاپ شده را به عنوان خسارت به دانشگاه تربیت مدرس، تادیه کند.

ماده ۵ : دانشجو تعهد و قبول می کند در صورت خودداری از پرداخت های بهای خسارت، دانشگاه مذکور را از طریق مراجع قضایی مطالبه و وصول کند، به علاوه به دانشگاه حق می دهد به منظور استیفای حقوق خود، از طریق دادگاه، معادل وجه مذکور در ماده ۴ را از محل توقیف کتابهای عرضه شده نگارنده برای فروش، تامین نماید.

ماده ۶ : اینجانب مرضیه ابراهیمی دانشجوی رشته آمار زیستی مقطع کارشناسی ارشد تعهد فوق و ضمانت اجرایی آن را قبول کرده، به آن ملتزم می شوم.

نام و نام خانوادگی
تاریخ و امضا



۹۴، ۳، ۱۵

تقدیم

بر پدر و مادر و خواهر عزیزم

آنان که باور ناب بودن، لذت زیستن، جسارت خواستن، عظمت رسیدن و تمام تجربهای یکتای زندگی، مدیون حضور سبزشان است.

مشکرو قدردانی

شکرشایان نثار ایزدمنان که توفیق را رفیق را هم نمود و به همنشینی رحروان علم و دانش مقننم ساخت. اکنون که پایان نامه حاضر به انجام رسیده است، بر خود لازم می دانم مراتب سپاس را از بزرگوارانی به جا آورم که اگر دست یاریکشان نبود، هرگز این تحقیق با انجام نمی رسید. از اولین و بهترین معلمان زندگیم پدر و مادر مهربانم که بیماری پر مهرشان دشواری راه را برایم هموار ساخت. بسی شایسته است قدردانی کنم از استاد که تقدیر جناب آقای دکتر انوشیروان کاظم نژاد که زحمت راهنمایی این پایان نامه را بر عهده داشتند و همچنین از جناب آقای دکتر غلامی که همواره با راهنمایی های عالمانه و ارزشمند، مرا مورد لطف و محبت قرار داده اند. از اساتید گرامی جناب آقای دکتر اسحق و جناب آقای دکتر زایری که زحمت داوری این پایان نامه را منتقل شدند، کمال سپاس را دارم.

با احترام

مرضیه ابراهیمی

چکیده

زمینه و اهداف: در تحقیقات پزشکی مدلسازی و پیش‌بینی از اهمیت زیادی برخوردار است. روش‌های پیش‌بینی را می‌توان با تکنیک‌هایی بهبود داد. مدل‌های آماری و شبکه عصبی مصنوعی از مدل‌هایی هستند که در رده‌بندی و پیش‌بینی مورد استفاده قرار می‌گیرند، اما مدل‌های آماری نیاز به پیش‌فرض‌هایی دارند و شبکه عصبی نیازمند حجم نمونه کافی برای آموزش است. از این رو در این پایان‌نامه اختلاط خبره‌ها را معرفی می‌کنیم که یکی از مدل‌های رایج در ترکیب دسته‌بندها است. ترکیب دسته‌بندها روشی برای بهبود کارایی در مسائل رده‌بندی است که دارای تعداد الگوی محدودتری هستند. این مدل‌ها برای رگرسیون و رده‌بندی مفیدند، ولی ما تنها به دسته‌بندی توسط این مدل پرداخته‌ایم. در اختلاط خبره‌ها فضا بین چند دسته‌بند تقسیم می‌شود و با یک شبکه میانجی اجرا می‌گردد. اختلاط خبره‌ها با استفاده از روش‌های تکراری نظیر الگوریتم ماکسیمم مورد انتظار و شبکه‌های عصبی قابل‌برازش است.

روش بررسی: این مطالعه شامل ۲۱۳ بیمار هیپاتیت سی (۱۹۵ مرد، ۱۸ زن، فاصله سنی ۱۲-۶۶) است. متغیر پاسخ در این مدل بهبودی پس از درمان در نظر گرفته شده است. در ابتدا مدل لجستیک به داده‌ها برازش داده شد و فاکتورهای موثر شناخته شده‌اند و بعد از آن مدل‌های شبکه عصبی و اختلاط خبره‌ها برازش داده شد. برای مقایسه عملکرد این سه مدل از منحنی مشخصه عملکرد و صحت پیش‌بینی استفاده شده است.

یافته‌ها: نتایج حاصل از مدل لجستیک نشان‌دهنده معناداری متغیرهای شاخص توده بدنی، سن، بارویروسی، RS12 است. در تحلیل این داده‌ها با دو مدل لجستیک و شبکه عصبی، سطح پیشگویی مشابهی را نشان داده است. مدل اختلاط خبره‌ها توسط شبکه عصبی برازش داده شد. این مدل در سطح عملکرد پیشگویی همانند شبکه عصبی بوده است. شبیه‌سازی‌ها این مطلب را تایید کردند که مدل اختلاط خبره‌ها یک مدل وابسته به ساختار داده است.

نتیجه: اختلاط خبره‌ها در مقایسه با یک مدل رده‌بندی عملکرد را بهبود می‌بخشد. اما این مدل به ساختار داده‌ها وابسته است.

کلید واژه: رده‌بندی، شبکه عصبی مصنوعی، اختلاط خبره‌ها، هیپاتیت C.

| | |
|---|----|
| فصل ۱ مقدمه و مروری بر مطالب گذشته | ۱ |
| ۱-۱ مقدمه | ۲ |
| ۲-۱ اهداف پایان نامه | ۴ |
| ۳-۱ نگاهی گذرا به پایان نامه | ۴ |
| ۴-۱ هوش مصنوعی | ۵ |
| ۵-۱ یادگیری ماشین | ۶ |
| ۶-۱ پیش پردازش داده‌ها | ۸ |
| ۷-۱ بیان مسئله | ۱۰ |
| ۸-۱ داده‌های مورد بررسی در پایان نامه | ۱۰ |
| فصل ۲ مواد و روش‌ها | ۱۴ |
| ۱-۲ مقدمه | ۱۵ |
| ۲-۱-۲ رگرسیون لجستیک | ۱۷ |
| ۳-۱-۲ شبکه عصبی | ۱۸ |
| ۲-۲ تحلیل نتایج در مدل‌های رده‌بندی | ۲۵ |
| ۱-۲-۲ روش اعتبارسنجی متقابل | ۲۵ |
| ۲-۲-۲ صحت رده‌بندی و خطاهای رده‌بندی | ۲۵ |
| ۳-۲-۲ انواع خروجی دسته‌بندها | ۲۷ |
| ۳-۲ سیستم‌های دسته‌بند چندگانه | ۲۹ |
| ۱-۳-۲ دلایل استفاده از سیستم‌های دسته‌بند چندگانه | ۲۹ |
| ۲-۳-۲ خروجی در ترکیب چند دسته‌بند پایه | ۳۱ |
| ۴-۲ اختلاط خبره‌ها | ۳۶ |
| ۵-۲ الگوریتم ماکسیمم مورد انتظار | ۳۸ |
| ۶-۲ مدل‌بندی اختلاط خبره‌ها با شبکه عصبی | ۴۳ |
| ۱-۶-۲ الگوریتم آموزش اختلاط خبرگان | ۴۴ |
| ۳-۶-۲ اختلاط خبرگان رقابتی و تعاونی | ۴۷ |
| ۷-۲ روش‌های ایجاد گوناگونی در دسته‌بندهای پایه | ۴۸ |
| ۱-۷-۲ روش‌های ضمنی | ۴۸ |
| ۲-۷-۲ روش‌های صریح | ۴۹ |
| ۳-۷-۲ معیارهای کمی برای محاسبه گوناگونی | ۵۰ |
| ۸-۲ شبیه‌سازی داده‌ها | ۵۲ |
| فصل ۳ یافته‌ها | ۵۴ |
| ۱-۳ پردازش داده‌ها | ۵۵ |

| | |
|---------|--|
| ۵۹..... | ۲-۳ مدل لجستیک..... |
| ۶۱..... | ۳-۳ شبکه عصبی..... |
| ۶۵..... | ۴-۳ مقایسه شبکه عصبی و مدل لجستیک..... |
| ۶۸..... | ۵-۳ اختلاط خبره‌ها..... |
| ۶۸..... | ۶-۳ نتایج شبیه‌سازی داده‌ها..... |
| ۷۱..... | فصل ۴ بحث، نتیجه گیری و پیشنهادات..... |
| ۷۶..... | فهرست منابع..... |
| ۷۸..... | ضمائم..... |
| ۸۵..... | چکیده انگلیسی..... |

| | |
|---|----|
| شکل ۱-۱ روند بازشناسی الگو..... | ۶ |
| شکل ۱-۲ شبکه عصبی پرسپترون یک لایه | ۲۱ |
| شکل ۲-۲ استدلال آماری برای ترکیب دسته بندها. * D بهترین دسته‌بند در این مسئله، منحنی یرونی ناحیه‌ای که همه دسته‌بندها را شامل می‌شود و فضای مشخص شده ناحیه ای است که دسته بندهای مناسب در آن قرار می‌گیرد. | ۲۸ |
| شکل ۳-۲ خط‌چین‌ها مسیری که هریک از دسته‌بندها در هنگام آموزش طی می‌کند را نمایش می‌دهد. | ۲۹ |
| شکل ۴-۲ تقسیم بندی روشهای ترکیب از دیدگاه آموزش پذیری ترکیب کننده..... | ۳۱ |
| شکل ۵-۲ شمای از یک یادگیرنده متا | ۳۵ |
| شکل ۶-۲ شمایی از مدل اختلاط خبره‌ها | ۳۴ |
| شکل ۷-۲ شمایی از یک شبکه عصبی با یک لایه پنهان | ۳۸ |
| شکل ۸-۲ نمودار آموزش خبره‌ها و شبکه میانجی با استفاده از تابع خطا..... | ۳۴ |
| شکل ۹-۲ شبیه‌سازی جامعه ای باخوشه‌هایی نرمال..... | ۵۳ |
| شکل ۱۰-۲ شبیه‌سازی جامعه متداخل از ۳ خبره | ۶۶ |
| شکل ۱۳-۳ نمودار راک مقایسه دو مدل لجستیک اختلاط خبره‌ها..... | ۵۳ |
| شکل ۱-۳ درصد مشاهدات جهش ژنی در دو گروه بیماران | ۵۴ |
| شکل ۲-۳ درصد مشاهدات سطح بیماری بین دو گروه بیماران | ۵۴ |
| شکل ۳-۳ درصد مشاهدات جنس در دو گروه | ۵۵ |
| شکل ۴-۳ منحنی راک مدل لجستیک | ۵۷ |
| شکل ۵-۳ تغییرات حساسیت شبکه در تغییر تعداد گره لایه میانی | ۵۸ |
| شکل ۶-۳ تغییرات صحت پیش‌بینی در نمونه‌های آموزش و آزمون | ۵۹ |
| شکل ۷-۳ نمودارهای تغییرات صحت پیش‌بینی در نرخ یادگیری‌های مختلف | ۶۰ |
| شکل ۸-۳ نمودارهای تغییرات صحت پیش‌بینی در مقابل تغییرات مومنتوم در شبکه عصبی MLP .. | ۶۴ |
| شکل ۹-۳ تغییرات صحت پیش‌بینی در تعداد گره‌های متفاوت | ۶۲ |
| شکل ۱۰-۳ منحنی راک برای دو مدل شبکه عصبی و رگرسیون لجستیک | ۶۴ |

| | | |
|-----------|---|----|
| جدول ۱-۲ | تعداد تصمیم های مشابه بین دودسته بند | ۵۱ |
| جدول ۳-۱۱ | آماره های توصیفی متغیرهای کمی | ۵۶ |
| جدول ۲-۳ | توزیع فراوانی متغیرهای کیفی در دو گروه | ۵۶ |
| جدول ۳-۳ | برآورد ضرایب مدل لجستیک مرتبط با بهبودی بیماران | ۶۰ |
| جدول ۳-۴ | نتایج اجرای شبکه در ۱۰ مرتبه تکرار اعتبارسنجی | ۶۵ |
| جدول ۳-۵ | صحت کلاس بندی مدل شبکه عصبی و مدل لجستیک در پیش بینی بهبودی بیماران | ۶۶ |
| جدول ۳-۶ | مقایسه فاصله اطمینان ۹۵٪ دو مدل شبکه عصبی مصنوعی و رگرسیون لجستیک | ۶۶ |
| جدول ۳-۷ | فاکتورهای موثر بر بهبود بیماران بر اساس دو مدل لجستیک و شبکه عصبی | ۶۷ |
| جدول ۳-۸ | صحت کلاس بندی مدل شبکه عصبی و مدل اختلاط خبره ها در پیش بینی بهبودی بیماران | ۶۹ |
| جدول ۳-۹ | صحت پیش بینی دربرازش مدل های مختلف | ۶۹ |
| جدول ۳-۱۰ | توزیع نمونه ها در مراحل یادگیری هر خبره | ۷۰ |

فصل اول

مقدمه و مروری بر مطالعات گذشته

۱-۱ مقدمه

خلاصه‌ای از واقعیت را مدل گویند. به بیان دیگر، نمایش مجرد یا فیزیکی یک سیستم یا سامانه را از یک دیدگاه و نگاه خاص مدلسازی می‌نامند. مدل‌ها، انواع گوناگونی مانند مدل فیزیکی، مدل ریاضی، مدل آماری، مدل گرافی، نرم‌افزاری، و... دارند و کاربردهای حیاتی متنوع و فراوانی در همه زمینه‌های علوم و فناوری دارند. در این میان مدلسازی آماری نسبت به سایر انواع مدلسازی دارای تنوع بیشتری می‌باشد. روشهایی نظیر مدل‌های رگرسیونی (خطی و غیر خطی)، تحلیل خوشه‌ای، آنالیز تشخیصی، سری‌های زمانی و... از روشهای مدلسازی آماری هستند، در این میان سه عامل نوع متغیرهای توضیحی، توزیع احتمالی متغیر پاسخ و رابطه موجود میان متغیرهای کمکی با متغیر پاسخ از عوامل مهم برانتخاب نوع مدل می‌باشد. مثلاً وقتی با پاسخ‌های پیوسته مواجه هستیم رگرسیون خطی معمولی و یا مدل‌های تحلیل واریانس می‌توانند روش‌های مناسبی برای پیشبرد اهداف مطالعه باشند. حال هر چه متغیر پاسخ حساس تر باشد دقت پیش‌بینی صحیح نیز از اهمیت بیشتری برخوردار می‌شود. بخصوص در حیطه علوم پزشکی و بهداشتی که در آن سلامت افراد مدنظر محققان می‌باشد. از این رو محققان همواره دنبال روش‌های دقیقتر برای پیش‌بینی بوده‌اند که از آن جمله می‌توان به روش شبکه عصبی مصنوعی اشاره نمود.

شبکه‌های عصبی مصنوعی یا به زبان ساده‌تر شبکه‌های عصبی سیستم‌ها و روش‌های محاسباتی نوینی هستند برای یادگیری ماشینی، نمایش دانش، و در انتها اعمال دانش به دست آمده در جهت پیش‌بینی پاسخ‌های خروجی از سامانه‌های پیچیده است. ایده اصلی این گونه شبکه‌ها (تا حدودی) الهام‌گرفته از شیوه کارکرد سیستم عصبی زیستی، برای پردازش داده‌ها، و اطلاعات به منظور یادگیری و ایجاد دانش قرار دارد. عنصر کلیدی این ایده، ایجاد ساختارهایی جدید برای سامانه پردازش اطلاعات است. این سیستم از بی شمار عناصر پردازشی فوق‌العاده بهم‌پیوسته با نام نورون تشکیل شده که برای حل یک مسئله با هم هماهنگ عمل می‌کنند و توسط سیناپس‌ها (ارتباطات الکترومغناطیسی) اطلاعات را منتقل می‌کنند. در این شبکه‌ها اگر یک سلول آسیب ببیند بقیه سلول‌ها می‌توانند نبود آنرا جبران کرده، و نیز در بازسازی آن سهیم باشند. این شبکه‌ها قادر به یادگیری‌اند. یادگیری در این سیستم‌ها به صورت تطبیقی صورت می‌گیرد، یعنی با استفاده از مثال‌ها وزن سیناپس‌ها به گونه‌ای تغییر می‌کند که در صورت دادن ورودی‌های جدید، سیستم پاسخ درستی تولید کند. توافق دقیقی بر تعریف شبکه عصبی در میان محققان وجود ندارد؛ اما اغلب آنها موافقند

که شبکه عصبی شامل شبکه‌ای از عناصر پردازش ساده (نورونها) است، که می‌تواند رفتار پیچیده کلی تعیین شده‌ای از ارتباط بین عناصر پردازش و پارامترهای عنصر را نمایش دهد. مطالعات نشان داده اند که شبکه عصبی برای تخمین^۱ و تقریب^۲ کارایی بسیار بالایی دارند. در برآزش شبکه های عصبی مصنوعی از الگوریتم تکراری تا رسیدن به مینیمم خطا استفاده می‌شود. یکی از مشکلات استفاده از این الگوریتم گرفتار شدن در مینیمم های محلی تابع خطا است. بمنظور رفع این مشکل می‌توان از سیستم‌های شورایی استفاده کرد.

در جوامع برای تصمیم‌گیری‌های مهم در مواقع کلیدی مردم به گروهی که در مجلس شورا و یا هیاتی که مدیریت را بعهدہ گرفته اند اعتماد می‌کنند و این مجالس بر اساس این تفکر که دو فکر بهتر از یک فکر عمل می‌کند شکل می‌گیرند. در سال ۱۸۱۸ لاپلاس نشان می‌دهد که ترکیب مناسبی از دو روش احتمالاتی بهتر از مدلی با یک مؤلفه انجام می‌شود. روش هایی که از ترکیب چند دسته‌بند تشکیل می‌شوند به روش های شورایی معروف هستند [۱].

یکی از روش های جالب شورایی^۳، اختلاط خبره ها^۴ است که توان بالقوه در اجرای ماشین‌های یادگیری ایجاد می‌کند. این مدل بر اساس اصل تقسیم و تسخیر^۵ عمل می‌کنند [۲]. در این مدل فضای مسئله بین چند خبره^۶ با نظارت مدخل شبکه^۷ تقسیم می‌شود. این مدل در سال ۱۹۹۱ توسط جیکوبز و جردن معرفی شده است [۲].

اختلاط خبره‌ها ابتدا در رشته های کامپیوتر و الکترونیک به کار برده شد و سپس در سیستم های هوشمند^۸ و سیستم‌های تصمیم‌گیری^۹ در شبکه عصبی و تحلیل و پردازش تصاویر و اطلاعات کاربرد داشته است و از دیگر کاربردهای آن در علوم پزشکی به منظور تشخیص و پیش‌بینی میزان خطر بروز یک بیماری را می‌توان نام برد. از کاربردهای این روش در سال‌های اخیر در مطالعه در تشخیص بیماری دیابت و تشخیص سیروز کبدی و همچنین در مطالعاتی بر روی سرطان سینه قابل ذکر است [۵و۴]. هریس و همکارانش شبکه اختلاط خبره‌ها را برای رسیدن به کارایی بالاتر در مدل های رگرسیون تعمیم داده اند. آنان معیار توقف را رسیدن به کمترین خطای کلی در شبکه قرار داده اند [۶].

¹Estimation

²Approximation

³Ensamble

⁴Mixture of experts

⁵Divided and conquer

⁶Experts

⁷Getting network

⁸Intelligent system

⁹Decision support system

منظور از خبره آمیخته همان اختلاط خبره‌ها است.

۱-۲ اهداف پایان نامه

هدف از انجام این پایان نامه، به کارگیری مدل اختلاط خبره‌ها در پیش‌بینی داده‌ها است. این مدل، معمولاً جهت رده‌بندی داده‌ها با حجم کم بکار رفته و قابل ذکر است که صحت پیش‌بینی این مدل نسبت به یک مدل رده‌بندی افزایش چشمگیری را در مقالات نشان داده است و از این رو، اجرای آن بر روی داده‌های مورد نظر انجام گرفته است. در این مطالعه به بررسی پارامترهای مختلف تاثیرگذار بر مدل پرداخته شده است.

۱-۳ نگاهی گذرا به پایان نامه

در این فصل یادگیری ماشین و ضروریات در انجام این روش بررسی شده است و همچنین داده‌های مورد استفاده در پایان نامه معرفی می‌شوند.

فصل دوم به بیان روش‌های رده‌بندی داده‌ها پرداخته شده است و نحوه تحلیل نتایج در مدل‌های رده‌بندی و اعتبار سنجی دسته‌بندها^۱ معرفی می‌شوند، از روشهای متداول رده‌بندی می‌توان به آنالیز تشخیصی^۲ و مدل‌های لجستیک^۳ و شبکه‌های عصبی مصنوعی^۴ اشاره داشت که با توجه به نوع مسأله از توانمندی‌های متفاوتی برخوردارند. شبکه عصبی یکی از متداولترین ابزارهای دسته‌بندی و پیش‌بینی در یادگیری ماشین است. این روش به عنوان مدل اصلی در مطالعه حاضر، به طور مفصل پرداخته شده است. برآورد پارامترها در شبکه عصبی با استفاده از الگوریتم یادگیری تکراری محاسبه می‌شود، این الگوریتم در جستجوی کمترین خطای شبکه هستند که ممکن است در کمینه‌های محلی گرفتار شود. یادگیری شورایی^۵ روشی است که به منظور بهبود دقت مدل رده‌بندی استفاده می‌گردد. بنابراین در ادامه به معرفی سیستم‌های دسته‌بند چندگانه^۶ پرداخته، و روش‌های برآورد این مدل‌ها بیان شده است.

در اختلاط خبره‌ها، هر مدل دسته‌بندی پایه با توجه به مقدار پارامترهایش، به پاسخ متفاوتی برای مسأله می‌رسد و با ترکیب پاسخ‌ها، دقت مدل افزایش می‌یابد. به منظور ترکیب مناسبی از خروجی چند مدل، ابتدا باید در جستجوی چارچوب ریاضی برای ترکیب این مدل‌ها بود تا از نقاط قوت آنها استفاده شود و از نقاط ضعف آنها پرهیز گردد، همچنین طبق تحقیقات انجام گرفته، وجود دسته‌بند های پایه با کارایی مناسب و تقریباً مستقل از یکدیگر از ملزومات این مدل می‌باشد [۷].

در متون بررسی شده از روش‌های آموزش مختلفی در مدل اختلاط خبره‌ها بهره گرفته شده

¹Classifier

²Discriminant analysis

³Logestic

⁴Artificial neural network

⁵Ensamble learning

⁶Multiple classifier systems

است. از جمله روش‌های آموزش این سیستم‌ها، استفاده از الگوریتم ماکسیمم مورد انتظار^۱ می‌باشد که در بعضی مقالات از این الگوریتم برای برآورد پارامترهای مدل اختلاط خبره‌های سلسله‌مراتبی^۲ استفاده شده است [۸-۱۰]. از دیگر روش‌های برآورد مدل استفاده از شبکه عصبی چند لایه پرسپترون است که در بسیاری از مقالات از این روش استفاده شده است [۱۱]. همچنین استفاده از برآورد کننده ییزی در آموزش اختلاط خبره‌ها از دیگر روش‌های بکار برده شده است [۲].

۴-۱ هوش مصنوعی

هوش مصنوعی قابلیت‌ای است برای انجام فعالیت‌ها و وظایف رایانه که همانند انسان نیاز به دانش، دقت، استدلال، آموزش، فهم قابلیت‌های ادراکی دارد.

ریشه اصلی سیستم‌های خبره^۳ و یا سیستم‌های مبتنی بر دانش^۴ (KBS) به حوزه مطالعاتی به نام هوش مصنوعی^۵ (AI) برمی‌گردد. تحقیق درباره‌ی AI در دهه ۱۹۴۰ شروع شد، زمانی که اولین نسل رایانه‌ها در مؤسسات تحقیقاتی به کار برده شد، اساس منطق ماشین مبتنی بر علم ریاضیات بوده است و افرادی مانند گرت گودل، الونزو چاچ، آلون ترینگ با استفاده از روش‌های فرموله شده‌ای در استدلال منطقی به آن پرداخته بودند.

تحقیق اصلی این افراد بر محاسبه پیش‌بینی و قیاس منطقی^۶ متمرکز شده است. این نحوه فرموله کردن نقش مهمی در سیستم‌های AI داشته است از پیامدهای توجه به AI در حل مسائل از حوزه ریاضی به حوزه مسائل دنیایی واقعی با تمرکز به فنون حل مسائل کلی معمول از این طریق، ظهور KBS و سیستم‌های خبره تجاری بوده است.

سیستم‌های مبتنی بر دانش شامل عملگرهایی است که مشخص می‌کند که چطور یک سیستم از یک وضعیت می‌تواند به وضعیت بعد و نهایتاً به سوی وضعیت هدف پیش رود. برای مثال در یک بازی شطرنج عملگرها قواعدی هستند که مهره‌های مختلف را قادر به تغییر مکان از یک خانه به خانه دیگر می‌کند. از نقایص مدل‌های موجود این بود که با افزایش اندازه مسائل فضای جستجوی آنها به صورت نمایی افزایش می‌یافت، بنابراین محققین با تمرکز بر دانش خاص در هر مساله، توانستند فضای جستجو را محدود تر نمایند. در سال ۱۹۸۶ واترمن، برنامه‌ای با کیفیت بالا بنحوی که دانش خاصی در حوزه آن مساله مرتبط باشد، را طراحی کرد. اولین فعالیت‌های KBS در حوزه زندگی واقعی نظیر تشخیص عفونی و پیش‌بینی منابع طبیعی معدنی در نواحی مختلف جغرافیایی دنیا بوده است [۱۲].

¹Expectation – maximization algorithm

²Hierarchical mixtures of experts

³Expert systems

⁴Knowledge-based systems

⁵Artificial intelligence

⁶Logical reasoning

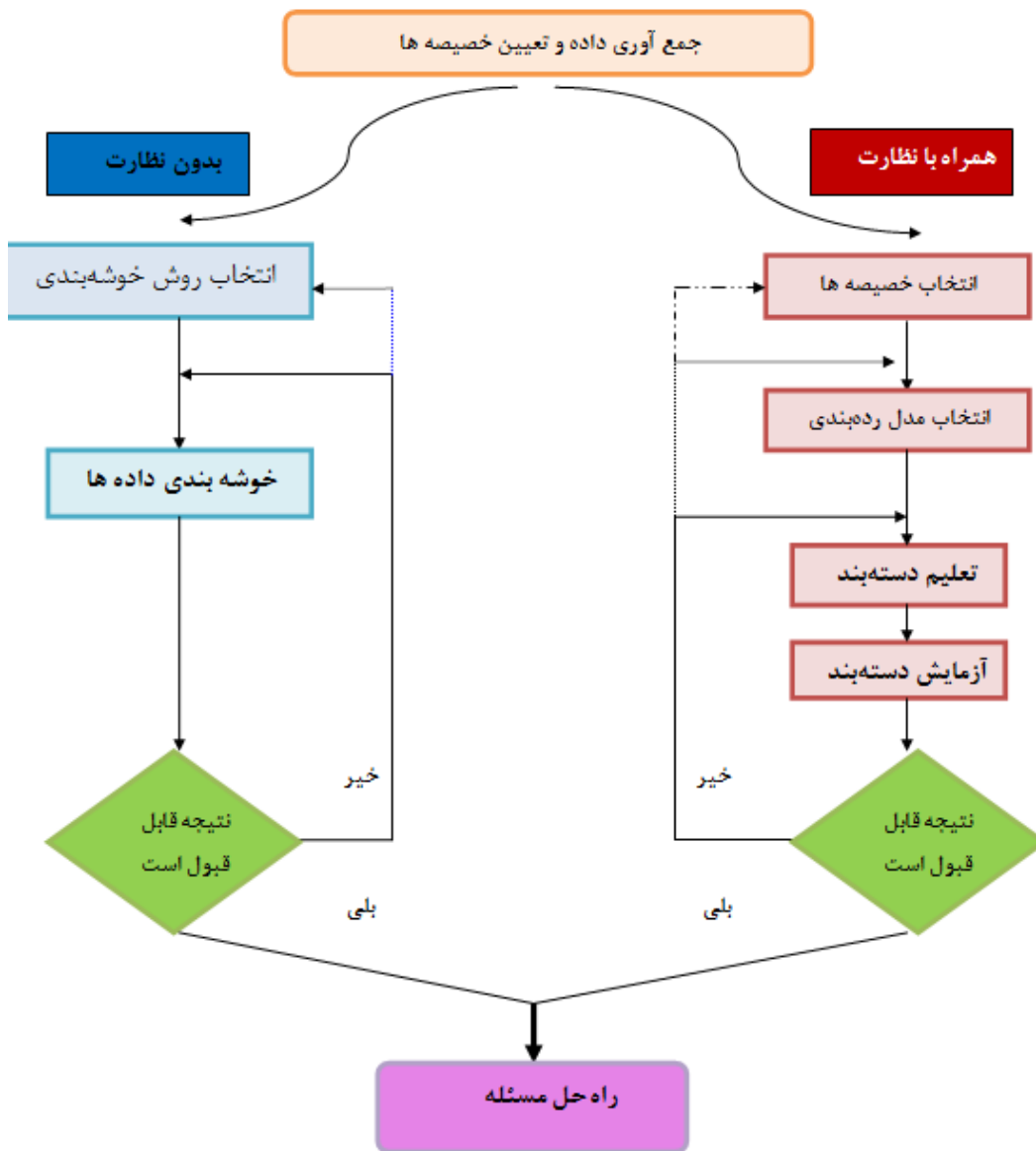
۱-۵ یادگیری ماشین

در سال ۱۹۵۹ اصطلاح یادگیری ماشین اولین بار توسط، ساموئل مطرح شد که به بررسی توانایی آموزش رایانه بدون وجود برنامه نویسی صریح می‌پردازد. به عبارت دیگر یادگیری ماشین راه هایی را که رایانه می‌تواند به طور مستقیم از داده ها دانش به دست‌آورد را بررسی می‌کند .

در سال ۱۹۶۳ مورگان و همکارانش در مقابل فرض های محدود کننده آمار کلاسیک از روش درخت رگرسیونی، که روشی بدون فرض های آماری می‌باشد را معرفی کردند . این روش به عنوان اولین روش غیر آماری یادگیری ماشین در حل مسائل آماری معروف می‌باشد. این الگوریتم ها با هدف بهبود عملکرد از طریق تجربه عمل می‌کنند. در حدود سالهای ۱۹۸۰ برهم کنش بین دو رشته آمار و یادگیری ماشین شروع شد، محققین یادگیری ماشین با سه مسأله رگرسیون، رده‌بندی، خوشه‌بندی توسط آماردانها آشنا شدند سپس آنها شروع به جستجو راه‌حل های غیر آماری و ناپارامتری جهت حل این مسائل کرده‌اند [۱۳].

یادگیری ماشین یک شاخه مهم از گرایش هوش مصنوعی است ، که ترکیبی از روش های ابتکاری هوش مصنوعی به همراه تحلیل آماری پیشرفته می‌باشد. که هدف از آن تعلیم یک ماشین است بطوریکه بتواند تجربیات و نمونه های موجود را یاد بگیرد. حاصل این یادگیری ایجاد یک مدل طبقه بندی است که بر اساس آن ماشین می‌تواند نمونه هایی را که در آینده می‌بیند و مشابه نمونه های موجود هستند در کلاس مناسب خود قرار دهد [۱۴]. امروزه روشهای بازشناسی الگو^۱ ، به عنوان یک شاخه از یادگیری ماشین، کاربردهای فراوانی در زمینه های مختلف علمی و صنعتی پیدا کرده‌اند. در فرآیند بازشناسی الگو، الگوهای ورودی در کلاسها و دسته‌های از پیش تعیین شده دسته‌بندی می‌شوند. روند بازشناسی الگو در شکل ۱-۱ نشان داده شده است [۱۵].

¹Pattern Recognition



شکل ۱۰-۱ روند بازشناسی الگوشکل

اولین گام در بازشناسی الگو، جمع آوری تعداد مناسبی نمونه از الگوهای مورد نظر است. این بخش زمان زیادی از فرآیند طراحی سیستم بازشناسی الگو را به خود اختصاص می‌دهد و گاهی اوقات با مشکلاتی همراه است. پس از جمع آوری نمونه‌های لازم، باید اقدام به انتخاب نوع ویژگی کرد. انتخاب نوع ویژگی نیازمند دانش اولیه در مورد داده‌ها است. انتخاب باید بگونه‌ای باشد که ویژگی‌ها، وجه مشترک الگوهای متعلق به یک کلاس و وجه تمایز الگوهای کلاسهای مختلف باشند. توانمندی ویژگی برای جداسازی نمونه‌های کلاسهای مختلف، معیار انتخاب آن است. مرحله بعدی انتخاب خصیصه^۱ها به معنیدست آوردن خصوصیات الگوهاست، به نحوی که این

^۱ Features

خصوصیات وجه مشترک الگوهای متعلق به یک کلاس و وجه تمایز الگوهای کلاس‌های مختلف باشند. پس از انتخاب مدل رده‌بندی، پارامترها باید مشخص شوند. روش‌های مختلفی برای تعیین این پارامترها ارائه شده‌اند. گروهی از این روشها بر اساس خواص آماری، مانند بردار میانگین و ماتریس کواریانس نمونه‌های آموزشی هر کلاس در هر دسته‌بندی، پارامترهای طبقه‌بندی را تعیین می‌کنند. در مقابل گروه دیگری از روشها، با استفاده از الگوریتم‌های یادگیری تکراری مانند آموزش یک طبقه بند شبکه عصبی بر روی نمونه‌های یادگیری، پارامترهای رده‌بندی را تعیین می‌کنند. در این روشها پارامترهای یادگیری در طول فرآیند یادگیری تعیین می‌شوند. در یادگیری بدون ناظر، پس از جمع آوری داده‌های کافی از الگوهای مورد نظر، یک روش خوشه‌بندی انتخاب می‌شود و بر اساس آن، داده‌ها خوشه‌بندی می‌شوند. اگر نتیجه خوشه‌بندی قابل قبول نباشد، با تغییر روش خوشه‌بندی یا تغییر پارامترهای خوشه‌بندی، کار ادامه پیدا می‌کند تا نتیجه مورد نظر حاصل شود.

یکی از ایده‌های جالب در یادگیری ماشین تجزیه یک مسئله مشکل به زیر مسئله‌هایی ساده‌تر و حل این مسائل با تکنیک‌های استاندارد است. به این سیستم هافرا تحلیل^۱ گفته می‌شود. یکی از زیر شاخه‌های این سیستم‌ها، سیستم‌های یادگیرنده چندگانه است. که در آن محققین از یادگیرنده‌های چندگانه برای یک مسئله یکسان، و ترکیب نتیجه آنها استفاده می‌کند. گروه یادگیرنده‌ها را در این مسائل شورا نامیده می‌شود. روشهای متعددی برای یادگیری ماشین و تشکیل شوراها و ترکیب تصمیم‌ها وجود دارد [۱۵].

در بسیاری از مسائل پیش‌بینی‌کننده احتیاج به تصمیم‌گیری در موقعیت‌های مختلف دارد مانند راننده‌ای که برای خیره شدن نیاز به یادگیری رانندگی در جاده، در بارندگی و جاده‌های یخ زده دارد. در مسائل رده‌بندی نیز هر یک از دسته‌بندها نقاط قوت و ضعف متفاوت دارند و رقابت و مدیریت این دسته‌بندها در سیستم‌های شورایی با توجه به درک این محدودیت‌ها در پیش‌بینی داده‌ها به دست می‌آید. محدودیت در اینجا به معنای ضعف در پیش‌بینی است که موجب خطا می‌شود.

۱-۶ پیش پردازش داده‌ها

کیفیت داده‌ها در استخراج نتایج مطلوب و اطلاعات حقیقی بسیار موثر است، پایگاه‌های داده اغلب دارای داده‌های پرت، گمشده و ناایستا هستند. بنابراین برای ارتقاء کیفیت داده‌ها لازم است در ابتدای کار داده‌ها به صورت زیر پردازش شوند.

پاکسازی داده‌ها: برای انجام یک تحلیل مطلوب لازم است مقادیر گمشده جایگزین شوند، و داده‌های مزاحم شناسایی و به نحوی مناسب با آنها برخورد و ناایستایی‌ها اصلاح شوند.

الف) مقادیر گمشده وجود مقادیر گمشده در داده‌ها می‌تواند تحلیل داده‌ها را بسیار دشوار سازد. در صورت وجود مقادیر گمشده در داده‌ها باید به گونه‌ای مناسب در مورد آنها تصمیم

¹Meta Analysis