

چکیده

زنجیره‌های تامین امروزی که پیوسته توسط جهانی سازی، برون سپاری، پیکربندی‌های محصول، نوسانات تقاضا، هزینه‌ها و SKUهای متنوع به چالش کشیده می‌شوند، فوق‌العاده پیچیده هستند. همه این فشارها شرکت‌ها را وادار می‌کند راهی برای ارضای نیازهای مشتریان بیابند که به موقع، کارآمد و سودآور باشد. گسترش سطح دسترسی به اینترنت و افزایش میل به خرید آنلاین، تعداد مشتریان را افزایش داده است. در شرایطی که تنوع کالا و تعداد مشتریان زیاد است، حل مسائلی مانند تحویل به موقع کالا یا خدمات، انتخاب و تعیین سفارش‌ها در انبارهای غیرمتمرکز، تخصیص انبار به مشتریان و... دشوار است. برای حل این چالش‌ها، استفاده از مدل سازی ریاضی با روش‌های حل فراابتکاری پیشنهاد شده است. اما به دلیل تعداد زیاد حالت‌های تخصیص، حل مدل‌های ریاضی بسیار پیچیده و زمان بر است. با بهبود قدرت محاسباتی و فضای ذخیره سازی، روش‌های مبتنی بر داده برای حل این چالش‌ها مورد مطالعه و بررسی قرار گرفته است. از این رو در این پژوهش نیز از روش داده کاوی برای تجزیه و تحلیل داده‌های مربوط به زنجیره تامین یک فروشگاه زنجیره‌ای با فروش حضوری و آنلاین استفاده شده است و نتایج به دست آمده نشان می‌دهد که الگوریتم‌های یادگیری ماشین رگرسیون و دسته بندی به ترتیب عملکرد فوق العاده‌ای در پیش بینی میزان فروش هر سفارش و وجود ریسک ارسال با تأخیر سفارش دارند. این موضوع نشان می‌دهد که استفاده از الگوریتم‌های یادگیری ماشین می‌تواند عملکرد بسیار مؤثری در حوزه زنجیره تامین فروشگاه داشته باشند.

واژگان کلیدی: عارضه یابی، زنجیره تامین، فروشگاه زنجیره‌ای، داده کاوی، یادگیری ماشین

فهرست مطالب

۱- مقدمه	۸
۲- پیشینه تحقیق	۱۰
۱-۲- مرور ادبیات	۱۰
۲-۲- جدول مرور ادبیات	۱۶
۳-۲- شکاف تحقیق	۲۳
۳- مطالعه موردی	۲۴
۴- مبانی نظری	۲۷
۱-۴- زنجیره تامین	۲۷
۲-۴- عارضه یابی	۲۷
۳-۴- مدیریت زنجیره تامین	۲۷
۴-۴- داده کاوی	۲۸
۱-۴-۴- انواع داده کاوی	۲۸
۵-۴- مدل CRISP-DM	۳۰
۱-۵-۴- مرحله اول: درک کسب و کار	۳۱
۲-۵-۴- مرحله دوم: بررسی و درک داده ها	۳۱
۳-۵-۴- مرحله سوم: آماده سازی یا پیش پردازش داده ها	۳۲
۴-۵-۴- مرحله چهارم: مدل سازی	۳۳
۵-۵-۴- مرحله پنجم: تست و ارزیابی مدل	۳۳
۶-۵-۴- مرحله ششم: توسعه مدل نهایی و استقرار	۳۳
۶-۴- یادگیری ماشین	۳۴
۱-۶-۴- انواع یادگیری ماشین	۳۴
۱-۱-۶-۴- الگوریتم های یادگیری ماشین با نظارت	۳۴
۲-۱-۶-۴- الگوریتم های یادگیری ماشین بدون نظارت	۳۵
۳-۱-۶-۴- یادگیری تقویتی	۳۵
۲-۶-۴- بیش برازش، کم برازش و برازش مناسب	۳۵
۳-۶-۴- موازنه واریانس و بایاس	۳۷
۵- روش تحقیق	۳۹
۱-۵- پیش پردازش داده ها: پاک سازی داده، کاهش ابعاد، انتخاب ویژگی ها	۳۹

- ۴۰-۵-۲- تجزیه و تحلیل اکتشافی داده‌ها..... ۴۰
- ۴۰-۵-۳- استخراج قوانین انجمنی..... ۴۰
- ۴۱-۵-۳-۱- الگوریتم Apriori..... ۴۱
- ۴۲-۵-۳-۲- معیارهای ارزیابی قوانین انجمن..... ۴۲
- ۴۲-۵-۴- پیاده‌سازی الگوریتم‌های یادگیری ماشین..... ۴۲
- ۴۲-۵-۴-۱- آماده‌سازی مجموعه داده برای پیاده‌سازی الگوریتم‌های یادگیری ماشین..... ۴۲
- ۴۴-۵-۴-۲- تقسیم مجموعه داده به سه قسمت آموزشی، اعتبارسنجی و تست..... ۴۴
- ۴۴-۵-۴-۳- پیاده‌سازی الگوریتم‌های یادگیری ماشین برای پیش‌بینی میزان فروش هر سفارش..... ۴۴
- ۴۴-۵-۴-۳-۱- رگرسیون خطی..... ۴۴
- ۴۶-۵-۴-۳-۲- رگرسیون ستیغی..... ۴۶
- ۴۷-۵-۴-۳-۳- رگرسیون لاسو..... ۴۷
- ۴۸-۵-۴-۳-۴- رگرسیون درخت تصمیم..... ۴۸
- ۴۸-۵-۴-۳-۵- رگرسیون جنگل تصادفی..... ۴۸
- ۴۹-۵-۴-۳-۶- رگرسیون بیزی..... ۴۹
- ۴۹-۵-۴-۴- پیاده‌سازی الگوریتم‌های یادگیری ماشین برای پیش‌بینی وجود ریسک ارسال با تاخیر... ۴۹
- ۴۹-۵-۴-۴-۱- درخت تصمیم‌گیری..... ۴۹
- ۵۱-۵-۴-۴-۲- جنگل تصادفی..... ۵۱
- ۵۲-۵-۴-۴-۳- دسته‌بندی کیسه‌ای..... ۵۲
- ۵۲-۵-۴-۴-۴- دسته‌بندی تقویتی گرادیان..... ۵۲
- ۵۹-۵-۴-۴-۵- دسته‌بندی XGboost..... ۵۹
- ۶۰-۵-۵- اعتبارسنجی و بررسی صحت الگوریتم‌های رگرسیون و دسته‌بندی..... ۶۰
- ۶۰-۵-۵-۱- میانگین مربعات خطا..... ۶۰
- ۶۰-۵-۵-۲- جذر میانگین مربعات خطا..... ۶۰
- ۶۱-۵-۵-۳- میانگین قدرمطلق خطا..... ۶۱
- ۶۱-۵-۵-۴- ضریب امتیاز R^2 ۶۱
- ۶۱-۵-۵-۵- مربع R تنظیم شده..... ۶۱
- ۶۲-۵-۵-۶- ماتریس اغتشاش..... ۶۲
- ۶۲-۵-۵-۷- دقت..... ۶۲
- ۶۳-۵-۵-۸- صحت..... ۶۳

۶۳	۵-۵-۹- پوشش
۶۳	۵-۵-۱۰- امتیاز F_1
۶۳	۵-۵-۱۱- اعتبارسنجی متقابل
۶۴	۵-۵-۱۲- منحنی ROC – AUC
۶۶	۵-۶- بررسی نتایج و مقایسه الگوریتم‌ها
۶۷	۶- یافته‌های تحقیق
۶۷	۶-۱- تجزیه و تحلیل اکتشافی
۷۶	۶-۲- استخراج قوانین انجمنی
۷۸	۶-۳- ارزیابی و مقایسه عملکرد الگوریتم‌های یادگیری ماشین
۸۰	۷- نتیجه‌گیری و پیشنهاد
۸۲	۸- منابع

فهرست اشکال

۲۸	شکل (۱-۴) انواع داده کاوی
۲۹	شکل (۲-۴) متدولوژی‌های داده کاوی
۲۹	شکل (۳-۴) متدولوژی‌های داده کاوی
۳۰	توسعه مدل نهایی و استقرار
۳۰	شکل (۴-۴) مدل داده کاوی CRISP-DM
۳۲	شکل (۵-۴) داده ساختاریافته و غیرساختاریافته
۳۲	شکل (۶-۴) پاک‌سازی داده
۳۳	شکل (۷-۴) اعتبارسنجی داده
۳۸	شکل (۸-۴) مصورسازی موازنه واریانس و بایاس برای برآوردگر
۳۹	شکل (۱-۵) روش تحقیق مطالعه
۴۵	شکل (۲-۵) مثالی از پیاده‌سازی روش رگرسیون خطی
۵۱	شکل (۳-۵) مثالی از پیاده‌سازی الگوریتم درخت تصمیم برای طبقه‌بندی داده‌ها
۵۱	شکل (۴-۵) مثالی از پیاده‌سازی الگوریتم جنگل تصادفی برای دسته‌بندی داده‌ها

شکل (۵-۵) نحوه عملکرد الگوریتم دسته‌بندی کیسه‌ای	۵۲
شکل (۶-۵) مثالی از مسئله دسته‌بندی دو کلاسه	۵۴
شکل (۷-۵) نمایش صفحه پیش‌بینی به شکل سه‌بعدی	۵۴
شکل (۸-۵) شکل نمایش باقی‌مانده‌ها	۵۵
شکل (۹-۵) درخت ساخته شده برای متغیرهای x و باقی‌مانده r	۵۵
شکل (۱۰-۵) نمایش پیش‌بینی بروز شده	۵۷
شکل (۱۱-۵) نمایش باقی‌مانده‌های بروز شده	۵۷
شکل (۱۲-۵) درخت ساخته شده برای متغیرهای x و باقی‌مانده r بروز شده	۵۸
شکل (۱۳-۵) نمایش پیش‌بینی بروز شده	۵۸
شکل (۱۴-۵) فرایند الگوریتم تقویتی گرادیان	۵۹
شکل (۱۵-۵) نحوه بهینه‌سازی در الگوریتم XGboost	۵۹
شکل (۱۶-۵) نحوه عملکرد اعتبارسنجی متقابل مونت کارلو	۶۴
شکل (۱-۶) همبستگی ویژگی‌ها در مجموعه داده	۶۷

فهرست جداول

جدول (۱-۳) تعریف ویژگی‌های مجموعه داده پاک‌سازی شده	۲۴
جدول (۱-۵) فرایند کدگذاری وان‌هات	۴۳
جدول (۲-۵) ماتریس اغتشاش	۶۲
جدول (۱-۶) ترتیب قوانین انجمنی استخراج شده بر اساس معیار پشتیبانی از زیاد به کم	۷۷
جدول (۲-۶) قوانین انجمنی استخراج شده از مجموعه داده، ترتیب بر اساس معیار Lift از زیاد به کم	۷۸
جدول (۳-۶) مقایسه الگوریتم‌های رگرسیون یادگیری ماشین برای پیش‌بینی فروش هر سفارش	۷۸
جدول (۴-۶) مقایسه الگوریتم‌های دسته‌بندی یادگیری ماشین برای پیش‌بینی وجود ریسک ارسال با تأخیر	۷۹

فهرست نمودارها

نمودار (۱-۴) منحنی بیش برآزش بر اساس چندجمله‌ای مرتبه ۱۵.....	۳۶
نمودار (۲-۴) منحنی کم برآزش بر اساس چندجمله‌ای مرتبه ۱.....	۳۶
نمودار (۳-۴) منحنی برآزش مناسب بر اساس چندجمله‌ای مرتبه ۴.....	۳۷
نمودار (۱-۵) منحنی مشخصه عملکرد برای سه روش مختلف دسته‌بندی.....	۶۵
نمودار (۲-۵): نواحی مطلوب و نامطلوب در منحنی ROC.....	۶۵
نمودار (۱-۶) متغیر هدف وجود ریسک ارسال با تأخیر.....	۶۸
نمودار (۲-۶) نمودار انواع معاملات انجام شده.....	۶۸
نمودار (۳-۶) نمودار روزهای ارسال واقعی محصول خریداری شده.....	۶۹
نمودار (۴-۶) نمودار وضعیت تحویل سفارش‌ها در مجموعه داده.....	۷۰
نمودار (۵-۶) نمودار فراوانی دسته‌بندی محصولات سفارش داده شده.....	۷۱
نمودار (۶-۶) نمودار بخش‌بندی مشتریان فروشگاه.....	۷۲
نمودار (۷-۶) نمودار فراوانی دپارتمان‌های فروشگاه.....	۷۲
نمودار (۸-۶) نمودار فراوانی بازار محل تحویل سفارش.....	۷۳
نمودار (۹-۶) نمودار فراوانی حالت‌های حمل‌ونقل سفارش‌ها.....	۷۴
نمودار (۱۰-۶) نمودار توزیع سفارش‌ها فروشگاه طی روزهای هفته.....	۷۵
نمودار (۱۱-۶) نمودار حالت‌های حمل‌ونقل سفارش‌ها به تفکیک وجود ریسک ارسال با تأخیر.....	۷۶
نمودار (۱۲-۶) منحنی ROC برای الگوریتم‌های دسته‌بندی بکار گرفته شده.....	۷۹

۱- مقدمه

تشخیص صحیح و به موقع مسائل ریشه‌ای سازمان، اولین گام به منظور ایجاد تحول، حفظ بقا و برتری سازمان‌ها است. در بسیاری از مباحث مدیریتی، سازمان به بدن انسان تشبیه می‌شود، چرا که سازمان نیز یک موجود پویا و زنده است. همان‌طوری که پیش‌نیاز انجام هر درمان و بهبود در بدن، انجام آزمایش‌های کاملی از وضعیت بدن انسان است و هرچه دقت این آزمایش‌ها بیشتر و دقیق‌تر باشد بهبودها و فرایند درمان مؤثرتر خواهد بود، در مورد ایجاد اصلاحات و حرکت به سوی رشد و پیشرفت سازمان نیز انجام فرایند عارضه‌یابی همین حکم را دارد و با انجام این فرایند سعی در یافتن فرصت‌های بهبودی و تنگناهایی خواهد شد که ممکن است به عنوان مانع حرکت سازمان، شناسایی گردد و در جهت بهبود و رفع این موانع برنامه‌ریزی و اقدام مناسب صورت گیرد.

در واقع زنجیره‌تأمین یکی از حوزه‌های بسیار وسیع بوده که زیرمجموعه‌های مختلفی را در بر دارد. در زنجیره تأمین بحث کنترل زمان و هزینه از اهمیت کلیدی برخوردار است. تحویل محموله‌ها و کالا بر اساس زمان‌بندی مشخص شده، جهت حصول اطمینان از اینکه در نهایت محصولات در زمان مناسب به دست مشتری می‌رسند، مورد توجه است. همچنین مدیریت و کنترل هزینه‌های کل شبکه تأمین، امری ضروری است، چرا که در صورت عدم کنترل مناسب و اتلاف هزینه‌ها در طول زنجیره، این هزینه‌ها بر روی قیمت تمام‌شده محصول یا خدمت ارائه شده اثر گذاشته و منجر به افزایش آن می‌گردد. در نهایت این افزایش قیمت تمام‌شده ممکن است به مشتری نهایی منتقل شده که مطلوب نیست و می‌تواند منجر به ایجاد نارضایتی و کاهش میزان تقاضا از جانب مشتریان و در نهایت کاهش فروش شرکت گردد. در حالت دیگر در صورت عدم افزایش قیمت محصول، با توجه به اتلاف‌های صورت گرفته و افزایش قیمت تمام‌شده، حاشیه سود شرکت کاهش پیدا کرده و این موضوع منجر به نارضایتی سهام‌داران می‌شود. از این رو در این بین داده‌کاوی می‌تواند به عنوان یک ابزار کمکی برای تحلیل بهتر هزینه‌ها و زمان‌ها در زنجیره‌تأمین و ارتقای کارایی آن مؤثر واقع شود. برخی از کاربردهای داده‌کاوی در فضای عارضه‌یابی زنجیره‌تأمین به صورت زیر می‌باشند:

- استفاده از تکنیک‌های داده‌کاوی برای شناسایی الگوهای مربوط تأخیر در تحویل محصولات یا ارائه خدمات

- استفاده از الگوریتم‌های داده‌کاوی به منظور کشف علل مربوط به افزایش هزینه‌های واقعی نسبت به هزینه‌های تخمینی

- و...

رویکرد داده‌محور، از مزیت دسترسی به منابع داده جدید برای کمک به شرکت‌ها در عارضه‌یابی زنجیره‌های تأمین استفاده می‌کند. زنجیره‌تأمین مدرن به طور فزاینده‌ای برای کسب و کارها، پیچیده و جهانی می‌شود. این بدان معناست که شناسایی آسیب‌پذیری‌های بالقوه در زنجیره‌تأمین، رفع تنگناها و واکنش سریع به اختلالات زنجیره تأمین، بسیار حیاتی است. زنجیره تأمین برای کسب و کارها، تنها یک زنجیره ساده

و خطی از فعالیت‌ها نیست. بلکه یک شبکه پویا از فرایندها، فناوری و افراد یکپارچه است. بسیاری از مدیران، توجه به زنجیره‌تأمین را تا زمانی که مشکل بزرگی پیش نیاید در نظر نمی‌گیرند، درحالی‌که سازمان‌ها باید استراتژی‌های منسجمی را توسعه داده و تصمیمات مبتنی بر داده را برای بهبود عملکرد زنجیره‌تأمین اتخاذ کنند. نقطه شروع برای بررسی زنجیره‌تأمین، تأیید این مسئله است که چگونه عملکرد زنجیره‌تأمین فعلی منجر به نتایج مالی و عملکردی بزرگ می‌شود. اتخاذ رویکردی سیستماتیک و ساختاریافته که معیارهای کلیدی عملکرد عملیاتی را با بازگشت سرمایه همسو می‌کند، اساسی است و می‌تواند نشان دهد که چگونه عملیات روزانه کسب‌وکار، منجر به نتایج مالی می‌شود. سپس با بینش مبتنی بر واقعیت، یک برنامه تحول متناسب که سودمندترین اقدامات را مورد هدف قرار می‌دهد، می‌تواند اتخاذ شود. بدون شک زنجیره تأمین داده‌محور می‌تواند برای بسیاری از شرکت‌ها فرصت رشد فراهم کند؛ چرا که اکثر شرکت‌ها، زنجیره‌های تأمین خود را تا زمانی که مشکلی پیش نیاید بررسی نمی‌کنند.

از رو در این تحقیق، یک راه‌حل مبتنی بر داده که از داده‌کاوی استفاده می‌کند؛ پیشنهاد شده است که بر روی داده‌های «یک فروشگاه زنجیره‌ای با فروش حضوری و آنلاین» پیاده‌سازی می‌شود؛ تا بررسی کند که آیا می‌توان مسائل و چالش‌های موجود در زنجیره‌تأمین را با داده‌کاوی حل نمود؟ یا خیر. همچنین تلاش می‌شود در ابتدا به بررسی و تجزیه و تحلیل اکتشافی داده‌ها پرداخته شود و در ادامه با استفاده از الگوریتم‌های یادگیری ماشین عوامل مؤثر بر فروش و سایر ویژگی‌های تأثیرگذار بر عملکرد زنجیره‌تأمین فروشگاه مدل‌سازی شود تا بتوان میزان فروش و همچنین عوامل مؤثر را پیش‌بینی و بررسی نمود.

۲- پیشینه تحقیق

۲-۱- مرور ادبیات

این مقاله به بررسی زنجیره تأمین خدمات آنلاین مبتنی بر داده از طریق دو دیدگاه دامنه و تأمین تقاضا می‌پردازد؛ همچنین عوامل مختلفی را که بر تقاضای خدمات آنلاین تأثیر می‌گذارند و راهبردهایی را که برای برآورده این تقاضا می‌توان در زنجیره تأمین به کار گرفت، مورد بررسی قرار داده و ویژگی‌های خدمات آنلاین و تأثیر آن‌ها بر زنجیره تأمین مانند سفارشی‌سازی خدمات و کنترل کیفیت خدمات را، تحلیل می‌کند؛ همچنین نقش داده‌ها در مدیریت زنجیره تأمین شامل جمع‌آوری، تجزیه و تحلیل و استفاده از داده‌ها، بررسی می‌شود. مقاله نشان می‌دهد که داده‌ها می‌توانند برای امور مختلفی چون بهینه‌سازی زنجیره تأمین و بهبود کیفیت خدمات به کار گرفته شوند. نویسندگان چالش‌ها و فرصت‌های زنجیره تأمین خدمات آنلاین مبتنی بر داده را نیز مورد بررسی قرار داده‌اند و همچنین، یک چارچوب برای مدیریت زنجیره تأمین خدمات آنلاین ارائه می‌دهند که شامل شناسایی عوامل تقاضا، توسعه راهبردهای تأمین تقاضا در زنجیره تأمین و بهره‌گیری از داده‌ها برای بهینه‌سازی زنجیره تأمین هستند [۱].

این مقاله یک رویکرد مبتنی بر داده برای همگام‌سازی تطبیقی تقاضا و عرضه در زنجیره‌های عرضه خرده‌فروشی omni-channel را پیشنهاد می‌کند. نویسندگان چالش‌هایی را که خرده‌فروش‌های omni-channel در برآوردن تقاضای مشتری با آن مواجه هستند مورد بحث قرار می‌دهند و استدلال می‌کنند که رویکردهای داده‌محور می‌توانند برای غلبه بر این چالش‌ها و بهبود عملکرد زنجیره تأمین مورد استفاده قرار گیرند. این مقاله چارچوبی برای همگام‌سازی تطبیقی ارائه می‌کند که شامل جمع‌آوری، تجزیه و تحلیل و تصمیم‌گیری در زمان واقعی است. نویسندگان استفاده از تکنیک‌های یادگیری ماشین را برای پیش‌بینی تقاضا و بهینه‌سازی تخصیص موجودی پیشنهاد می‌کنند. این مقاله شامل یک مطالعه موردی از یک خرده‌فروش برزیلی است که رویکرد پیشنهادی را اجرا کرد که منجر به بهبود دید زنجیره تأمین، کاهش انبارها و افزایش فروش شد. نویسندگان همچنین محدودیت‌ها و جهت‌گیری‌های تحقیقاتی آینده این رویکرد را مورد بحث قرار می‌دهند [۲].

این مقاله اندازه‌گیری، کاهش و پیشگیری از ضایعات مواد غذایی در زنجیره تأمین در خرید آنلاین را مورد بحث قرار می‌دهد. نویسندگان چالش‌های ضایعات مواد غذایی را در زنجیره‌های تأمین، از جمله اثرات اقتصادی، زیست‌محیطی و اجتماعی را برجسته می‌کنند و استدلال می‌کنند که خرید آنلاین می‌تواند به عنوان ابزاری برای کاهش ضایعات مواد غذایی با ارائه پیش‌بینی دقیق‌تر تقاضا، کاهش بیش از حد سفارش، و امکان مدیریت بهتر موجودی استفاده شود. این مقاله چارچوبی برای اندازه‌گیری ضایعات مواد غذایی در زنجیره‌های تأمین ارائه می‌کند که شامل کمی کردن مقدار ضایعات، شناسایی علل آن و ارزیابی اثربخشی استراتژی‌های کاهش است. نویسندگان استفاده از تجزیه و تحلیل داده‌ها و یادگیری ماشین را برای بهبود پیش‌بینی تقاضا و بهینه‌سازی مدیریت موجودی پیشنهاد می‌کنند. همچنین این پژوهش، شامل یک مطالعه

موردی از یک خواربارفروشی آنلاین است که رویکرد پیشنهادی را پیاده‌سازی کرده و منجر به کاهش ضایعات مواد غذایی و بهبود کارایی زنجیره تأمین شده است [۳].

این مقاله در مورد استفاده از تکنیک‌های یادگیری ماشین برای پیش‌بینی فروش در فروشگاه‌های خرده‌فروشی صحبت می‌کند و اشاره می‌کند که پیش‌بینی دقیق فروش برای مدیریت موجودی، برنامه‌ریزی عملیاتی و بهینه‌سازی درآمد ضروری است. آنها استفاده از الگوریتم‌های یادگیری ماشین مانند Random Forest، Support Vector Machine و Artificial Neural Network را برای پیش‌بینی فروش بر اساس داده‌های تاریخی و متغیرهای دیگر مانند تخفیف‌ها، تعطیلات و هواشناسی پیشنهاد می‌دهند. نویسندگان مزایا و محدودیت هر تکنیک را بررسی کرده و عملکرد آنها را با استفاده از معیارهایی مانند میانگین خطای درصدی مطلق مقایسه می‌کنند. آنها همچنین اهمیت انتخاب ویژگی و پیش‌پردازش داده را در بهبود دقت پیش‌بینی مورد بحث قرار می‌دهند. مقاله شامل یک مطالعه موردی از یک فروشگاه در هند است که رویکرد پیشنهادی پیاده‌سازی شده است و منجر به بهبود دقت پیش‌بینی فروش و مدیریت بهتر موجودی شده است [۴].

این مقاله به بررسی پیش‌بینی فروش خرده‌فروشان در بازارهای مد می‌پردازد و بیان می‌کند که پیش‌بینی دقیق فروش برای مدیریت مؤثر موجودی، برنامه‌ریزی تولید و رضایت مشتری حیاتی است. آنها چالش‌ها و فرصت‌های پیش‌بینی فروش خرده‌فروشان صنعت مد، از جمله چرخه‌های عمر کوتاه محصول، تغییرات سریع در اولویت‌های مصرف‌کننده، و نیاز به تجزیه و تحلیل داده‌های بلادرنگ را مورد بحث قرار می‌دهند. این مقاله مروری بر تکنیک‌های مختلف پیش‌بینی مورد استفاده در خرده‌فروشی مد، مانند روش‌های آماری، هوش مصنوعی و یادگیری ماشینی ارائه می‌کند و مزایا و محدودیت‌های هر روش را مورد بحث قرار داده و عملکرد آنها را با استفاده از معیارهایی مانند دقت پیش‌بینی و میانگین درصد خطا مقایسه می‌کند. آنها همچنین در مورد اهمیت کیفیت داده‌ها و پیش‌پردازش داده‌ها در بهبود دقت پیش‌بینی بحث می‌کنند. نویسندگان به این نتیجه رسیدند که هیچ راه‌حلی برای پیش‌بینی فروش خرده‌فروشی مد وجود ندارد و ممکن است ترکیبی از تکنیک‌های مختلف برای دستیابی به پیش‌بینی‌های دقیق و به موقع مورد نیاز باشد [۵].

این مقاله یک مطالعه موردی از مدیریت اختلالات در زنجیره تأمین خرده‌فروشی را ارائه می‌کند. نویسندگان اشاره می‌کنند که اختلالات زنجیره تأمین می‌تواند منجر به ضررهای مالی قابل توجه و آسیب به شهرت برند شود. آنها چالش‌های موجود در مدیریت اختلالات، از جمله نیاز به اطلاعات به موقع و دقیق، ارتباطات مؤثر و برنامه‌ریزی اضطراری را مورد بحث قرار می‌دهند. این مقاله یک مطالعه موردی از یک زنجیره تأمین خرده‌فروشی را توصیف می‌کند که در آن یک تأمین‌کننده اصلی آتش‌سوزی در تأسیسات تولیدی خود را تجربه کرد که باعث اختلال در عرضه یک دسته محصول حیاتی شد. نویسندگان گام‌های برداشته شده توسط شرکت خرده‌فروشی برای مدیریت اختلال را مورد بحث قرار می‌دهند. این مقاله همچنین استفاده از مدل‌های شبیه‌سازی زنجیره تأمین را برای ارزیابی سناریوهای مختلف و ارزیابی تأثیر اختلالات احتمالی

مورد بحث قرار می‌دهد. نویسندگان بر اهمیت مدیریت ریسک فعال و انعطاف‌پذیری زنجیره تأمین در به حداقل رساندن تأثیر اختلالات تأکید می‌کنند [۶].

این مقاله کاربرد یادگیری ماشین و مدل‌های ترکیبی را در پیش‌بینی فروش خرده‌فروشی بررسی می‌کند. نویسندگان اشاره می‌کنند که پیش‌بینی دقیق فروش برای خرده‌فروشان برای تصمیم‌گیری آگاهانه در مورد مدیریت موجودی و برنامه‌ریزی تولید بسیار مهم است. این مقاله مروری بر مدل‌های مختلف یادگیری ماشین، از جمله شبکه‌های عصبی، درخت‌های تصمیم‌گیری و رگرسیون برداری پشتیبان و مزایا و محدودیت‌های آن‌ها در پیش‌بینی فروش خرده‌فروشی ارائه می‌کند. نویسندگان همچنین یک مدل ترکیبی پیشنهاد می‌کنند که نقاط قوت مدل‌های مختلف را برای بهبود دقت پیش‌بینی ترکیب می‌کند و در مورد اهمیت کیفیت داده‌ها و پیش‌پردازش در آموزش مدل یادگیری ماشین بحث می‌کنند و بر نیاز به نظارت مداوم و به‌روزرسانی مدل برای در نظر گرفتن تغییرات در رفتار مصرف‌کننده و روند بازار تأکید می‌کنند. این مقاله نتیجه‌گیری می‌کند که یادگیری ماشین و مدل‌های ترکیبی پتانسیل بالایی در پیش‌بینی فروش خرده‌فروشی دارند و تحقیقات بیشتری برای کشف قابلیت‌های کامل آن‌ها در این زمینه مورد نیاز است [۷].

در این مقاله یک مدل بهینه‌سازی برای تخصیص منابع ارائه داده می‌شود و سپس یک الگوریتم مورچه^۱ بهبودیافته برای حل آن ایجاد خواهد شد. الگوریتم مورچه دارای مزایای همگرایی سریع به راه‌حل بهینه است و الگوریتم مورچه بهبودیافته نیز دارای مزایای آشکاری در تعادل انعطاف‌پذیری بهینه‌سازی‌های چندهدفه، تنظیم سرعت همگرایی و تنظیم پارامترهای عملیات می‌باشد.

در ادامه، اثربخشی و امکان‌سنجی روش و الگوریتم بهینه‌سازی با یک شبیه‌سازی عددی نشان داده می‌شود. مدل بهینه‌سازی نه تنها هزینه خدمات و زمان تحویل مورد نیاز را در تابع هدف منعکس می‌کند، بلکه بهینه‌سازی اثر اندازه خدمات و روابط مزایا و خطرات یکپارچه‌سازی را نیز در نظر می‌گیرد و نشان می‌دهد موفقیت خدمات خرید آنلاین مستلزم دو رابطه ضروری است؛ ۱- الگوهای مختلف خدمات سفارشی باید با طرح‌های مختلف تخصیص منابع زنجیره تأمین مطابقت داشته باشد و ۲- الگوهای مختلف ترکیب خدمات سفارشی شده نیز باید با طرح‌های تخصیص منابع زنجیره‌تأمین مختلف مطابقت داشته باشند [۸].

شرکت‌های خرید آنلاین B²C باید یکپارچه‌سازی منابع زنجیره‌تأمین^۲ خود را برای ارائه ظرفیت‌های خدمات، بهبود تجربیات مشتریان خود و معرفی خدمات سفارشی رضایت‌بخش بهینه کنند. این مطالعه با تجزیه و تحلیل ویژگی‌ها و حالت‌های خدمات SCRI در B²C، تعادل پویا بین ظرفیت‌های خدمات عرضه و تقاضا را در یک شرکت خرید آنلاین مورد بحث قرار می‌دهد. بحث در مورد این ظرفیت‌ها نه تنها باید اهداف

^۱ Ant algorithm

^۲ SCRI: supply chain resources integration

بهینه‌سازی سنتی را در نظر بگیرد، بلکه باید مناسب‌بودن منابع را از زوایای خدمات عمومی، خدمات اضطراری و ظرفیت‌های بالقوه استراتژیک برای دستیابی به یک SCRI مؤثر را ارزیابی کند؛ بنابراین، این مطالعه ثبات جهت‌گیری هدف ظرفیت منبع را به‌عنوان یک هدف مهم بهینه‌سازی در نظر گرفته و با شناسایی عوامل مشخصه ظرفیت و معرفی آن‌ها به SCRI، این ثبات را ارزیابی می‌کند. بدین ترتیب یک مدل بهینه‌سازی و یک الگوریتم مورچه بهبودیافته برای حل فرایند SCRI پیشنهاد شدند که امکان‌سنجی و اعتبار این مدل و الگوریتم با نشان‌دادن کاربرد آن‌ها و همچنین نشان‌دادن ارزش تحقیق آن‌ها تأیید شدند [۹].

این مقاله، به پیش‌بینی فروش با استفاده از مدل‌های رگرسیون لجستیک و K-نزدیک‌ترین همسایه^۳ بر اساس یادگیری ماشین می‌پردازد. اطلاعات از منابع داده‌های بین‌المللی جمع‌آوری شده است. پیش‌بینی فروش در توسعه کسب‌وکار، برآورد نیاز مالی، صحت تصمیمات مدیریتی، همکاری و هماهنگی، کنترل و غیره کاربرد دارد و این نشان از اهمیت پیش‌بینی فروش در کسب‌وکار دارد. پس از ایجاد و تست مدل‌ها بر روی داده‌ها، از ماتریس سردرگمی و روش اقلیدسی نتایج مورد ارزیابی قرار گرفتند. نتایج نشان می‌دهند که دو مدل رگرسیون لجستیک و K-نزدیک‌ترین همسایه می‌توانند عملکرد خوبی در پیش‌بینی فروش به‌ویژه در بازارهای محصولات FMCG که تقاضای بالایی دارند و به‌سرعت فروخته می‌شوند، داشته باشند [۱۰].

برای شرکت‌های خرده‌فروشی مدرن که زنجیره عظیمی از کسب‌وکارها را اداره می‌کنند، پیش‌بینی دقیق فروش، کلید توسعه شرکت‌ها و حتی موفقیت یا شکست آن‌ها است. پیش‌بینی فروش به شرکت‌ها این امکان را می‌دهد تا به طور مؤثر منابعی از جمله جریان نقدی، تولید و برنامه کسب‌وکار را تخصیص دهند. در این مقاله، یک مدل پیش‌بینی فروش کارآمد و دقیق با استفاده از یادگیری ماشین پیشنهاد می‌شود. در ابتدا، مهندسی ویژگی برای استخراج ویژگی‌ها از داده‌های فروش تاریخی انجام می‌شود. علاوه بر این، از XGBoost که یک مدل درختی تقویت‌کننده گرادیان توزیع‌شده بهینه‌سازی شده می‌باشد، برای پیش‌بینی میزان فروش آینده استفاده شده است. نتایج از طریق معیار خطای جذر میانگین مربعات^۴ که نوعی از معیار میانگین خطای مقیاس مطلق است، ارزیابی شده است. نتایج آزمایش بر روی مجموعه داده کالاهای خرده‌فروشی عمومی والمارت^۵ نشان می‌دهد که مدل پیشنهادی برای پیش‌بینی فروش با زمان محاسباتی و منابع حافظه کمتر بسیار خوب عمل می‌کند [۱۱].

پیش‌بینی فروش چالش‌برانگیزترین کار برای مدیریت موجودی، بازاریابی، خدمات مشتری و برنامه‌ریزی مالی کسب‌وکار صنعت خرده‌فروشی است. در این مقاله، یک تجزیه و تحلیل پیش‌بینی‌کننده از خرده‌فروشی مجموعه داده Citadel POS با استفاده از تکنیک‌های مختلف یادگیری ماشین انجام شده است. در این

^۳ K-Nearest Neighbour

^۴ RMSE: Root Mean Squared Error

^۵ Walmart

مطالعه، مدل‌های رگرسیون مختلف (رگرسیون خطی، رگرسیون جنگل تصادفی، رگرسیون افزایش گرادیان) و مدل‌های سری زمانی (ARIMA و LSTM)، برای پیش‌بینی فروش پیاده‌سازی شده است و تجزیه و تحلیل و ارزیابی پیش‌بینی دقیق ارائه گردیده است. مجموعه داده مورد استفاده در این تحقیق از Citadel POS سال ۲۰۱۳ تا ۲۰۱۸ به دست آمده است؛ نتایج نشان می‌دهد که XGboost از سری‌های زمانی و سایر مدل‌های رگرسیون بهتر عمل کرده و بهترین عملکرد را با میانگین خطای مطلق^۶ ۰.۵۱۶ و خطای ریشه میانگین مربع^۷ ۰.۶۳ به دست آورده است [۱۲].

رویدادهای اخیر، همانند همه‌گیری مداوم، عدم قطعیت زنجیره تامین را افزایش داده است. این عوامل ذاتی (مثل کمبود کانتینر) و عوامل خارجی (مثل افزایش تقاضا) عدم قطعیت‌ها در مدیریت زنجیره تامین را تشدید کرده و اهمیت BDA^۸ را در SCM^۹ بیشتر کرده است. بر این اساس، این مطالعه مروری سیستماتیک از مطالعات موجود در BDA را انجام می‌دهد. در حال حاضر، همراه با توسعه اخیر در زیرساخت‌های یادگیری ماشین و محاسبات، BDA در زنجیره تامین اهمیت زیادی پیدا کرده است. این مقاله چارچوبی از یک بررسی نظام‌مند ادبیات را از دیدگاه‌های بین‌رشته‌ای ارائه می‌کند. از دیدگاه سازمانی، این مطالعه به بررسی مبانی نظری و مدل‌های تحقیقاتی می‌پردازد که پایداری و عملکرد به دست آمده از طریق استفاده از تجزیه و تحلیل داده‌های بزرگ را توضیح می‌دهند؛ سپس، از دیدگاه فنی، این پژوهش انواع تجزیه و تحلیل داده‌های بزرگ، تکنیک‌ها، الگوریتم‌ها و ویژگی‌های توسعه یافته برای توابع زنجیره تامین افزایش یافته را تحلیل می‌کند و در نهایت، شکاف تحقیق را شناسایی کرده و جهت‌گیری‌های تحقیقاتی آتی را پیشنهاد می‌کند [۱۳].

این تحقیق تلاش می‌کند تا عوامل مؤثر بر عملکرد زنجیره تامین را برای تولیدکنندگان کوچک تحت سناریوی خرید الکترونیکی شناسایی کند. داده‌های اولیه مانند تقاضای مشتری از طریق توابع کنترل شده با ورودی تولید شده است و مفروضات بر اساس تجربه محقق تنظیم شده است. ادبیات موجود نشان می‌دهد که زمان کوتاه‌تر منجر به عملکرد بهتر زنجیره تامین می‌شود. با این حال، از طریق شبیه‌سازی پیشنهاد می‌شود که زمان طولانی‌تر، سطح خدمات را با افزایش موجودی ایمنی و سطح موجودی عمومی بهبود می‌بخشد. نتایج نشان‌دهنده وجود اثر شلاقی در زنجیره تامین خرید الکترونیکی به دلیل سفارش‌های عقب‌افتاده و تصمیمات مدیریتی است. مقایسه بین دو استراتژی کنترل موجودی مختلف با الگوهای تکمیل نشان می‌دهد که روش سطح سهام هدف در سناریوی خرید الکترونیکی نسبت به روش کمیت سفارش اقتصادی بهتر عمل می‌کند [۱۴].

در این مطالعه، یک مدل پیش‌بینی مبتنی بر خوشه‌بندی با ترکیب روش‌های خوشه‌بندی و یادگیری ماشین برای پیش‌بینی فروش خرده‌فروشی رایانه‌ای پیشنهاد شده است. ابتدا از تکنیک خوشه‌بندی برای

^۶ MAE: Mean Absolute Error

^۷ RMSE: Root Mean Square Error

^۸ Big Data Analytics

^۹ Supply Chain Management

خوشه‌بندی داده‌ها با ویژگی‌ها یا الگوهای مشابه در یک گروه استفاده شده است. همچنین، از تکنیک‌های یادگیری ماشین برای آموزش مدل پیش‌بینی هر گروه استفاده شده است. در مجموع شش مدل پیش‌بینی مبتنی بر خوشه‌بندی پیشنهاد شده است. داده‌های فروش واقعی برای رایانه‌های شخصی، رایانه‌های نوت‌بوک و نمایشگرهای کریستال مایع به‌عنوان مثال‌های تجربی استفاده شده‌اند. لازم به ذکر است که دقت پیش‌بینی را می‌توان با استفاده از مدل پیش‌بینی مبتنی بر خوشه‌بندی پیشنهادی افزایش داد. نتایج تجربی نشان می‌دهد که مدل ترکیبی از GHSOM و ELM عملکرد پیش‌بینی برتری را برای هر سه محصول در مقایسه با سایر روش‌ها داشته‌اند. این می‌تواند به طور مؤثر به‌عنوان یک مدل پیش‌بینی فروش مبتنی بر خوشه‌بندی برای خرده‌فروشی کامپیوتر استفاده شود [۱۵].

از آنجایی که صنعت مد به راحتی قابل پیش‌بینی نیست، پیش‌بینی دقیق فروش نیز ساده نیست. در این مطالعه، از روش‌های رگرسیون در یادگیری ماشین و تکنیک‌های تحلیل سری‌های زمانی برای پیش‌بینی میزان فروش بر اساس چندین ویژگی استفاده شده است. مدل‌ها بر روی داده‌های فروش والمارت در پلتفرم Microsoft Azure Machine Learning Studio اعمال شده‌اند. تکنیک‌های رگرسیون شامل رگرسیون خطی، رگرسیون بیزی، رگرسیون شبکه عصبی، رگرسیون جنگل تصادفی و رگرسیون درخت تصمیم تقویت شده بکار گرفته شده‌اند. علاوه بر این تکنیک‌های رگرسیون، روش‌های تحلیل سری‌های زمانی شامل ARIMA فصلی، ARIMA غیر فصلی، ETS فصلی، ETS غیر فصلی، روش متوسط و روش دریافت اجرا شده‌اند. نتایج می‌دهد که رگرسیون درخت تصمیم تقویت شده با ضریب تعیین ۰.۹۷ بهترین عملکرد را در این داده‌های فروش ارائه می‌دهد و عملکرد بهتری را در مقایسه با سایر روش‌ها برای پیش‌بینی فروش دارد [۱۶].

کلید موفقیت در تجارت، امروزه کنترل زنجیره تامین خرده‌فروشی است. پیش‌بینی تقاضای مشتری برای مدیریت زنجیره‌تأمین بسیار ضروری است. این مقاله، یک روش جدید با استفاده از یادگیری ماشین ارائه می‌کند که به پیش‌بینی دقیق تقاضا کمک می‌کند. در این مطالعه، داده‌های گذشته یک فروشگاه را جمع‌آوری شده است. از الگوریتم‌های K-Nearest Neighbor، Support Machine Vector، Gaussian، Decision Tree Classifier، Random Forest، Nave Bayes و رگرسیون‌ها در این مطالعه مقایسه شده‌اند. این مقاله، از ترکیب موقعیت فروشگاه، ماه و مناسبت آن ماه و سایر داده‌های مرتبط استفاده کرده است. مدل ارائه شده تقاضای آزمایشی برای یک محصول خاص را پیش‌بینی می‌کند. پس از ایجاد یک مجموعه داده و اعمال الگوریتم‌ها، نتایج و دقت الگوریتم‌های مختلف مقایسه شده‌اند. نتایج نشان می‌دهد که Gaussian Nave Bayes بهترین دقت را در بین الگوریتم‌ها دارد [۱۷].

این مقاله، یک چارچوب مفهومی برای درک نقشی که هوش مصنوعی می‌تواند در زنجیره ارزش خرده‌فروشی بازی کند، معرفی می‌کند. هدف اصلی این مقاله، درک بهتر یک زنجیره ارزش خرده‌فروشی مبتنی بر هوش مصنوعی است. به‌عنوان نقطه شروع، یک مرور مختصر از زنجیره ارزش سنتی خرده‌فروشی و فعالیت‌ها، ذی‌نفعان و فناوری درگیر در هر مرحله ارائه داده می‌شود. سپس ضعف‌های حاضر صنعت خرده‌فروشی توضیح داده می‌شود و به دنبال آن تمرکز ویژه‌ای بر نقشی که هوش مصنوعی در رفع اختلال

در این صنعت ایفا کرده است، مورد بررسی قرار می‌گیرد. سپس، فناوری‌های مختلف هوش مصنوعی مبتنی بر گارتتر را برای هر مرحله در زنجیره ارزش ترسیم می‌شود و نشان داده می‌شود که برخی از سرمایه‌گذاری‌های فناوری هوش مصنوعی می‌توانند اهداف متعددی را در زنجیره ارزش انجام دهند [۱۸]. در این مطالعه، یک تابع تقاضای زنجیره تامین خرید آنلاین بر اساس درجه همجوشی آنلاین به آفلاین با در نظر گرفتن آگاهی کم‌کربن خریداران آنلاین، با توجه به تصمیم‌گیری کم‌کربن زنجیره تامین، در ترکیب با پس‌زمینه عملیات آنلاین به آفلاین پیشنهاد شده است. بر این اساس، یک مدل تصمیم‌گیری کم‌کربن مبتنی بر فروشگاه‌های آنلاین و تامین‌کنندگان آن‌ها ایجاد شده و کاربرد مدل مورد تجزیه و تحلیل قرار می‌گیرد. نتایج نشان می‌دهد که ارتقای عملکرد کم‌کربن زمانی کارآمد است که ثبات عملیات آنلاین و آفلاین فروشگاه‌های آنلاین و تامین‌کنندگان آن‌ها بالا باشد. علاوه بر این، تصمیم‌گیری با در نظر گرفتن هزینه انتشار اطلاعات کربن برای بهبود آگاهی خریداران آنلاین نسبت به کربن پایین بسیار مفید است. اگر سطح کربن پایین کالاهای خرید آنلاین کمتر از استاندارد کم‌کربن باشد، بر این اساس، عملکرد کم‌کربن زنجیره تامین خرید آنلاین ترویج خواهد شد [۱۹].

این مقاله به بررسی اهمیت لجستیک در تجارت الکترونیک و نقش آن در مدیریت زنجیره تامین پرداخته است. همچنین بر روی نیاز به عملیات لجستیکی کارآمد و مؤثر برای اطمینان از رضایت مشتریان، بهبود فروش و حفظ مزیت رقابتی تأکید می‌کند. علاوه بر آن چالش‌هایی که صاحبان کسب‌وکار در لجستیک تجارت الکترونیک و مدیریت زنجیره تامین با آن روبرو هستند. نویسندگان این مقاله علاوه بر آنکه چندین راهبرد برای بهینه‌سازی لجستیک تجارت الکترونیک از جمله خودکارسازی، همکاری و سفارشی‌سازی ارائه می‌دهد؛ بر اهمیت تجزیه و تحلیل داده‌ها، ادغام فناوری و شیوه‌های پایدار در بهبود لجستیک و مدیریت زنجیره تامین تأکید می‌کنند و بر تأثیر لجستیک تجارت الکترونیک بر ذی‌نفعان مختلف، بهبود مداوم، نوآوری و سازگاری با تقاضاهای تغییرپذیر لجستیک تجارت الکترونیک و مدیریت زنجیره تامین تأکید می‌کند. در آخر؛ مقاله یک دیدگاه جامع از لجستیک تجارت الکترونیک در مدیریت زنجیره تامین، چالش‌های آن و راهکارهای پتانسیل آن را ارائه می‌دهد [۲۰].

این مقاله به بررسی رابطه بین کیفیت زنجیره تامین الکترونیکی و رضایت مشتری در خرید آنلاین می‌پردازد. نویسندگان یک مطالعه موردی بر روی یک وبسایت خرید آنلاین را انجام دادند و داده‌ها را با استفاده از پرسش‌نامه نظرسنجی جمع‌آوری کردند. آنها داده‌ها را با استفاده از روش‌های آماری تجزیه و تحلیل کردند و دریافتند که کیفیت زنجیره تامین الکترونیکی تأثیر مثبت قابل توجهی بر رضایت مشتری دارد. این مطالعه همچنین چندین عامل کلیدی که بر کیفیت زنجیره تامین الکترونیکی و رضایت مشتری در خرید آنلاین تأثیر می‌گذارد، مانند طراحی وبسایت، در دسترس بودن محصول، زمان تحویل، و خدمات مشتری را شناسایی می‌کند. نویسندگان توصیه‌هایی از جمله بهبود عملکرد وبسایت، بهینه‌سازی تدارکات و حمل و نقل و ارائه خدمات شخصی‌سازی شده را برای خرده‌فروشان آنلاین ارائه می‌کنند تا کیفیت زنجیره تامین

الکترونیکی و رضایت مشتری را بهبود ببخشند و بر اهمیت کیفیت زنجیره تأمین الکترونیکی در موفقیت خرده‌فروشان آنلاین و نیاز به بهبود مستمر برای برآورده کردن انتظارات مشتری تأکید می‌کند [۲۱].

این مقاله یک تحلیل پیش‌بینی فروش خرده‌فروشی با استفاده از تکنیک‌های یادگیری ماشینی ارائه می‌دهد. این مطالعه از داده‌های فروش تاریخی فروشگاه‌های در پاکستان برای توسعه مدل پیش‌بینی استفاده کرده است. نویسندگان پنج الگوریتم یادگیری ماشینی مختلف را به کار برده و عملکرد آن‌ها را مقایسه کردند. الگوریتم‌های استفاده شده عبارت‌اند از: Random Decision Tree، Multiple Linear Regression، Gradient Boosting، Forest و Artificial Neural Network.

از اندازه‌گیری‌های آماری مانند خطای مطلق میانگین درصد^{۱۰} و خطای میانگین مربعات ریشه^{۱۱} برای ارزیابی عملکرد مدل‌ها استفاده کردند. نتایج نشان داد که الگوریتم افزایش گرادیان با کمترین مقدار MAPE و RMSE بهترین عملکرد را داشت. نویسندگان معتقدند که مدل پیشنهادی می‌تواند توسط خرده‌فروشان برای پیش‌بینی فروش محصولات و تصمیم‌گیری‌های اطلاعاتی در رابطه مدیریت موجودی و استراتژی فروش استفاده شود [۲۲].

این مقاله یک چارچوب جدید برای پیش‌بینی فروش بر اساس مدل Convolutional Long Short-Term Memory ارائه می‌دهد. نویسندگان از مجموعه داده فروش روزانه یک فروشگاه زنجیره‌ای در چین برای آموزش و ارزیابی مدل استفاده کرده‌اند. آنها عملکرد این مدل را با مدل‌های سنتی پیش‌بینی سری زمانی مانند مدل ARIMA و ETS مقایسه کرده‌اند. نتایج نشان داد که مدل ConvLSTM در اصطلاح دقت و استحکام عملکرد مدل‌های سنتی را برتری داشته است. نویسندگان معتقدند که چارچوب پیشنهادی توسط خرده‌فروشان برای تصمیم‌گیری درباره مدیریت موجودی و استراتژی‌های فروش قابل استفاده است. این مطالعه اهمیت در نظر گرفتن عوامل خارجی مانند تعطیلات و رویدادهای تبلیغاتی در پیش‌بینی فروش را نیز بیان می‌کند. نویسندگان بیشتر می‌گویند که چارچوب می‌تواند به صنایع دیگر خارج از خرده‌فروشی مانند مالی و حمل‌ونقل توسعه یابد. به‌طور کلی، این مطالعه در مورد کاربرد تکنیک‌های یادگیری ماشین برای پیش‌بینی فروش در صنعت خرده‌فروشی ارزشمند است [۲۳].

این مقاله، مدل‌های پژوهشی‌ای ارائه داده است که می‌توانند با برخی تغییرات، برای اندازه‌گیری کارایی زنجیره تأمین خرده‌فروشی مورد استفاده قرار گیرند. مدل‌ها بر اساس رویکرد تحلیل پوششی داده‌ها DEA بودند. چهار گروه اصلی شناسایی شدند: مدل‌های استاندارد DEA، مدل‌های تجزیه کارایی، مدل‌های شبکه و مدل‌های مبتنی بر نظریه بازی. در بخش دوم مقاله، رویکردهای مختلف بر روی یک نمونه واقعی یک شرکت تجاری فعال در صربستان آزمایش شد. هفت زنجیره تأمین مشاهده شد که هر کدام از یک مرکز

^{۱۰} MAPE

^{۱۱} RMSE

توزیع DC و فروشگاه خرده‌فروشی RC تشکیل شده است. متغیرهای مورد استفاده عبارت بودند از تعداد مکان‌های پالت، هزینه‌های لجستیک، تعداد تحویل، دقت تحویل و گردش مالی، نتایج مزایا و معایب رویکردهای مختلف را در مثال واقعی نشان داد. سهم اصلی این مقاله در رویکردهای منحصربه‌فرد برای اندازه‌گیری کارایی زنجیره تامین خرده‌فروشی نهفته است [۲۴].

۲-۲- جدول مرور ادبیات

جدول (۲-۱) جدول مرور ادبیات

ردیف	عنوان مقاله	سال انتشار	نام مجله یا کنفرانس	نویسندگان	روش مورد استفاده	مطالعه موردی	یادگیری ماشین	الگوریتم										روش های ارزیابی			
								دسته بندی										Confusion Matrix	RMSE	MAE	RMSE
								linear regression	K-Nearest Neighbour	XGBoost	Random Forest	Decision Tree	SVM	سایر	خوشه بندی	خوشه بندی	سری های زمانی				
																	سایر				سایر
۱	Data-driven online service supply chain: a demand-side and supplyside perspective	۲۰۲۰	Journal of Enterprise Information Management	LeiLi et al.	بررسی زنجیره تامین از طریق دو دیدگاه دامنه تقاضا و تامین تقاضا																ارائه چارچوب برای مدیریت خدمات آنلاین
۲	A data-driven approach to adaptive synchronization of demand and supply in omni-channel retail supply chains	۲۰۲۰	International Journal of Information Management	Marina Meireles Pereira et al.	همگام سازی تطبیقی داده و تحلیل داده بلادرنگ		*														ارائه چارچوب برای همگام سازی تطبیقی

۳	Measurement, mitigation and prevention of food waste in supply chains: An online shopping perspective	۲۰۱۹	Industrial Marketing Management	Vasco Sanchez Rodrigues, et al.	یادگیری ماشین و تجزیه و تحلیل داده ها و کمی کردن داده های کیفی	*	*												الگوریتم طراحی شده به منظور شناسایی اتلاف مواد غذایی									ارائه روش ها برای بهبود پیش بینی تقاضا و بهینه سازی	
۴	Sales-forecasting of Retail Stores using Machine Learning Techniques	۲۰۱۸	International Conference on Computational Systems and Information Technology for Sustainable Solutions	Akshay Krishna et al.	پیش بینی فروش از طریق الگوریتم های یادگیری ماشین	*	*							*	*					مقایسه عملکرد روش های پیش بینی با استفاده از معیارها									
۵	A survey on retail sales forecasting and prediction in fashion markets	۲۰۱۵	Systems Science & Control Engineering	Samaneh Beheshti-Kashi et al.	پیش بینی بر اساس روش های مختلف	*													روش های آماری، هوش مصنوعی									مقایسه عملکرد راه حل های موجود	

[illegible]

روش اقلیدسی				*										*	*	*	*	پیش‌بینی فروش به کمک مدل‌های یادگیری ماشین	Neha Sehgal , Deepika Garg	International Journal of Innovative Research in Science , Engineering and Technology	۲۰۱۹	Sales Forecasting using Linear Regression and K-Nearest Neighbour	۱۰
			*									*				*	*	پیش‌بینی فروش به کمک مدل‌های یادگیری ماشین	Shilong, Zhang	۲۰۲۱ IEEE International Conference on Consumer Electronics and Computer Engineering (ICCECE)	۲۰۲۱	Machine learning model for sales forecasting by using XGBoost	۱۱
	*	*									*	*		*	*	*	*	پیش‌بینی فروش به کمک مدل‌های یادگیری ماشین	Sajawal, Muhamm ad et al.	Computer Science and Information Technology	۲۰۲۱	A Predictive Analysis of Retail Sales Forecasting using Machine Learning Techniques	۱۲
																		مرور ادبیات سیستماتیک از مطالعات موجود BDA در	Lee, In. Mangalar aj, George	Big Data and Cognitive Computing	۲۰۲۲	Big data analytics in supply chain management: a systematic literature review and research directions	۱۳
مقایسه آماری																*	*	شبیه‌سازی استراتژی	Li, Jiafu et al.	Computers &Industrial Engineering	۲۰۱۶	Impact of replenishment strategies on supply chain performance under e-shopping scenario	۱۴

۲-۳- شکاف تحقیق

باتوجه به بخش مرور ادبیات، نکته‌ای که حائز اهمیت است، حجم داده‌های مورد بررسی در مقالات است که غالباً از داده‌های محدودی استفاده شده است؛ بنابراین، استفاده از داده‌های حجیم‌تر برای افزایش دقت نتایج، به عنوان شکاف تحقیقاتی در این بررسی شناخته می‌شود. همچنین استفاده از الگوریتم‌های یادگیری ماشین به صورت ترکیبی تنها در یک مورد از مقالات، بررسی شده است و مطالعه مقالات نشان می‌دهد که همچنان نیاز به توسعه الگوریتم‌های یادگیری ماشین در این زمینه وجود دارد. به علاوه، تجزیه و تحلیل اکتشافی، در هیچ کدام از مقالات وجود ندارد و استخراج دانش به ندرت در مقاله‌های یادگیری ماشین صورت گرفته است.

در جدول مرور ادبیات مشاهده می‌شود که ۱۶ مورد از ۲۴ مقاله بررسی شده دارای مطالعه موردی هستند و باتوجه به اهمیت پیاده‌سازی روش‌ها در مطالعه موردی، در این پژوهش نیز از مطالعه موردی استفاده خواهد شد. براین اساس در این مطالعه، از مجموعه داده‌های یک فروشگاه زنجیره‌ای با فروش حضوری و اینترنتی استفاده شده است و تلاش می‌شود در ابتدا به بررسی و تجزیه و تحلیل اکتشافی داده‌ها پرداخته شود و در ادامه با استفاده از الگوریتم‌های یادگیری ماشین عوامل مؤثر بر فروش و سایر ویژگی‌های تأثیرگذار بر عملکرد زنجیره تامین فروشگاه مدل‌سازی شود تا بتوان میزان فروش و همچنین عوامل مؤثر را پیش‌بینی و بررسی نمود.

۳- مطالعه موردی

در این مطالعه، برای تجزیه و تحلیل زنجیره تامین یک فروشگاه با فروش حضوری و اینترنتی و پیاده سازی الگوریتم های یادگیری ماشین، ما از مجموعه داده اطلاعات زنجیره تامین یک فروشگاه با خرید اینترنتی و حضوری که توسط شرکت دیتاکو^{۱۲} بین سال های ۲۰۱۵ تا ۲۰۱۸، منتشر شده است، استفاده کرده ایم. این مجموعه داده در ابتدا شامل ۵۳ ستون (ویژگی^{۱۳}) و ۱۸۰۵۱۹ سطر (نمونه^{۱۴}) بوده است که پس از انتخاب ویژگی ها^{۱۵}، پاک سازی داده^{۱۶}، کاهش ابعاد^{۱۷} و استخراج ویژگی های مجموعه داده، تعداد ستون ها به ۶۰ ستون و سطرها به ۱۸۰۵۱۶ سطر رسید.

در ادامه، در جدول (۴-۱) ویژگی های مجموعه داده پاک سازی شده، به صورت مختصر توضیح داده شده است.

جدول (۳-۱) تعریف ویژگی های مجموعه داده پاک سازی شده

نام ویژگی	تعریف
Type	نوع معامله انجام شده
Days for shipping (real)	روزهای ارسال واقعی محصول خریداری شده
Days for shipment (scheduled)	روزهای تحویل برنامه ریزی شده محصول خریداری شده
Benefit per order	درآمد به ازای هر سفارش ثبت شده
Sales per customer	مجموع فروش انجام شده به ازای هر مشتری
Delivery Status	وضعیت تحویل سفارش ها: ارسال پیش از زمان، تأخیر در تحویل، لغو ارسال، ارسال به موقع
Late_delivery_risk	متغیر طبقه بندی که نشان می دهد اگر ارسال دیر باشد (۱) و اگر دیر نشده است (۰).
Category Id	کد دسته محصول
Category Name	نام دسته بندی محصول
Customer City	شهری که مشتری خرید را انجام داده است.
Customer Country	کشوری که مشتری خرید را انجام داده است.
Customer Id	شناسه مشتری
Customer Segment	بخش بندی مشتریان: مصرف کننده، شرکت های بزرگ، شرکت های کوچک
Customer State	ایالتی که فروشنده ای که خرید در آن ثبت شده است، به آن تعلق دارد.
Customer Zipcode	کد پستی مشتری

^{۱۲} Data-Co

^{۱۳} Feature

^{۱۴} Instance (Record)

^{۱۵} Feature Selection

^{۱۶} Data Cleaning

^{۱۷} Dimension Reduction

کد دپارتمان فروشگاه	Department Id
نام دپارتمان فروشگاه	Department Name
عرض جغرافیایی مربوط به محل فروشگاه	Latitude
طول جغرافیایی مربوط به محل فروشگاه	Longitude
بازار محل تحویل سفارش: آفریقا، اروپا، آمریکای لاتین، آسیا اقیانوسیه، ایالات متحده آمریکا	Market
شهر مقصد سفارش	Order City
کشور مقصد سفارش	Order Country
کد سفارش مشتری	Order Customer Id
تاریخ و ساعتی که در آن سفارش انجام شده است.	order date (DateOrders)
کد سفارش	Order Id
کد محصول تولید شده از طریق خواننده فرکانس رادیویی ^{۱۸}	Order Item Cardprod Id
مقدار تخفیف کالای سفارشی	Order Item Discount
درصد تخفیف کالای سفارشی	Order Item Discount Rate
کد کالای سفارشی	Order Item Id
قیمت محصولات بدون تخفیف	Order Item Product Price
نرخ سود کالای سفارشی	Order Item Profit Ratio
تعداد محصولات در هر سفارش	Order Item Quantity
میزان فروش	Sales
مبلغ کل به‌ازای هر سفارش	Order Item Total
سود سفارش در هر سفارش	Order Profit Per Order
منطقه‌ای از جهان که در آن سفارش تحویل داده می‌شود: آسیای جنوب شرقی، آسیای جنوبی، اقیانوسیه، آسیای شرقی، غرب آسیا، غرب ایالات متحده و غیره	Order Region
ایالتی منطقه‌ای که سفارش در آن تحویل داده می‌شود.	Order State
وضعیت سفارش: کامل، در انتظار، بسته، در انتظار پرداخت، لغو شده، در حال پردازش و غیره	Order Status
کد محصول	Product Card Id
کد دسته محصول	Product Category Id
نام محصول	Product Name
قیمت محصول	Product Price
تاریخ و زمان دقیق حمل‌ونقل	shipping date (DateOrders)
حالت‌های حمل‌ونقل شامل کلاس استاندارد، کلاس اول، کلاس دوم، همان روز	Shipping Mode
تاریخی که در آن سفارش انجام شده است.	DateOrders

^{۱۸} RFID: Radio Frequency Identification

زمانی که در آن سفارش انجام شده است.	TimeOrders
روزی از ماه که در آن سفارش انجام شده است.	DayOrders
ماهی که در آن سفارش انجام شده است.	MonthOrders
سالی که در آن سفارش انجام شده است.	YearOrders
روزی از هفته که در آن سفارش انجام شده است.	DayOfWeekOrders
ساعتی از روز که در آن سفارش انجام شده است.	HourOrders
دقیقه‌ای از ساعت که در آن سفارش انجام شده است.	MinutesOrders
تاریخ حمل‌ونقل	shipping Date
زمان حمل‌ونقل	shipping Time
روز حمل‌ونقل	shipping Day
ماه حمل‌ونقل	shipping Month
سال حمل‌ونقل	shipping Year
روز هفته حمل‌ونقل	DayOfWeek Shipping
ساعت حمل‌ونقل	shipping Hour
دقیقه‌ی حمل‌ونقل	shipping Minute

۴- مبانی نظری

۴-۱- زنجیره تامین

زنجیره تامین یک فرایند یکپارچه است که شامل تبدیل مواد خام به محصول نهایی و در نهایت تحویل به مشتری می‌شود. فرایند کامل زنجیره تامین به چهار سطح تقسیم می‌شود: تأمین کنندگان، تولید کنندگان، توزیع کنندگان و مشتریان (ماتور و همکاران ۲۰۱۷).

۴-۲- عارضه یابی

عارضه یابی مجموعه‌ای از فرایندها است که به مفهوم عامی اشاره می‌کند و آن چیزی نیست جز شناسایی و کشف معضلات و عارضه‌های سازمان در بخش‌های مختلف به منظور رفع و توسعه ظرفیت‌های سازمان. سازمان‌ها در گذر زمان با تأثیر گرفتن از تصمیمات مختلف و تغییرات محیطی درگیر مشکلات متعددی می‌شوند که برخی از این مشکلات در توسعه و رشد سازمان اثرگذار هستند. اغلب مشکلات و عارضه‌های سازمانی به صورت پنهان در سازمان باقی مانده و غالب آن‌ها به صورت محدودیت عمل می‌کنند. در واقع این مشکلات محدودیت نیستند؛ ولی چون رفع نشده‌اند به صورت محدودیت عمل کرده و فرایندهای سازمان را با اختلال همراه می‌کنند. حرکت سازمان به سمت بهبود و توسعه زمانی به طور صحیح و اثربخش انجام می‌گردد که عارضه‌های سازمانی، دلایل آن‌ها و اولویت دلایل به درستی شناسایی و مشخص گردند. از طرفی شناسایی نقاط ضعف یک سازمان زمانی می‌تواند برای سازمان اثربخش باشد که باعث شناسایی نقاط بهبود بالقوه گردد. پس از بررسی و شناسایی عارضه‌ها و نقاط ضعف سازمان، بهبودهای بالقوه شناسایی می‌شوند پس از این مرحله و با اولویت بندی، پروژه‌های بهبود برای بهبود در عملکردها و قابلیت‌های سازمان تعریف می‌گردد.

۴-۳- مدیریت زنجیره تامین

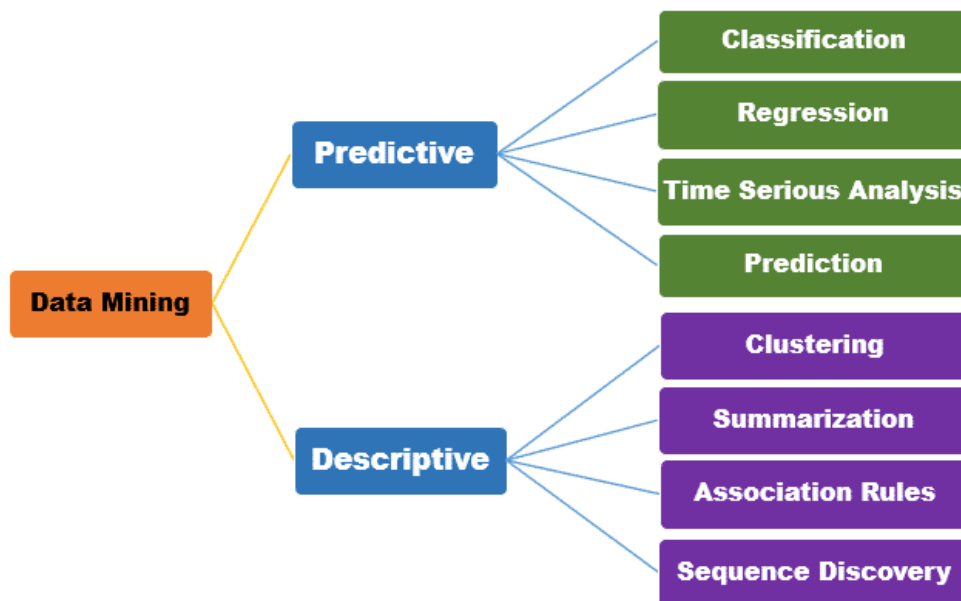
مدیریت زنجیره تامین، مدیریت جریان کالاها و خدمات، بین مشاغل و مکان‌ها است و شامل جابه‌جایی و ذخیره‌سازی مواد خام، موجودی در جریان کار و کالاهای نهایی و همچنین انجام سفارش از مبدأ تا نقطه مصرف می‌باشد. در واقع مدیریت زنجیره تامین ادغام فرایندهای تجاری کلیدی از کاربر نهایی از طریق تأمین کنندگان اصلی است که محصولات، خدمات و اطلاعاتی را ارائه می‌دهد که برای مشتریان و سایر ذی‌نفعان ارزش افزوده ایجاد می‌کند (لامبرت و همکاران، ۱۹۹۸؛ کوپر و همکاران، ۱۹۹۷). کریستوفر (۱۹۹۸) مدیریت زنجیره تامین را چنین تعریف کرد: «زنجیره تامین شبکه‌ای از سازمان‌ها است که از طریق پیوندهای بالادستی و پایین‌دستی درگیر فرایندهای مختلف و فعالیت‌هایی است که در قالب محصولات و خدمات در دست مشتری نهایی ارزش تولید می‌کنند.»

۴-۴- داده کاوی^{۱۹}

به مجموعه‌ای از روش‌های قابل‌اعمال بر پایگاه‌داده‌های بزرگ و پیچیده به‌منظور کشف الگوهای پنهان و جالب‌توجه نهفته در میان داده‌ها، داده کاوی گفته می‌شود. علم میان‌رشته‌ای داده کاوی، پیرامون ابزارها، متدولوژی‌ها و تئوری‌هایی است که برای آشکارسازی الگوهای موجود در داده‌ها مورد استفاده قرار می‌گیرند و گامی اساسی در راستای کشف دانش محسوب می‌شود.

داده کاوی که با عنوان «کشف دانش از داده»^{۲۰} نیز شناخته شده است، فرایند استخراج اطلاعات و دانش از داده‌های موجود در پایگاه‌داده یا انبار داده است.

۴-۱- انواع داده کاوی



شکل (۴-۱) انواع داده کاوی

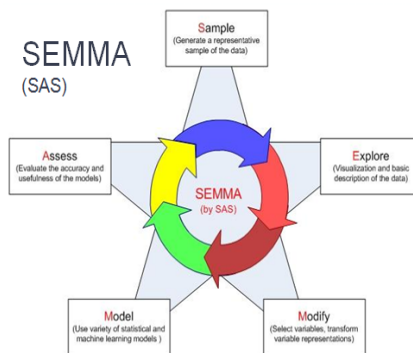
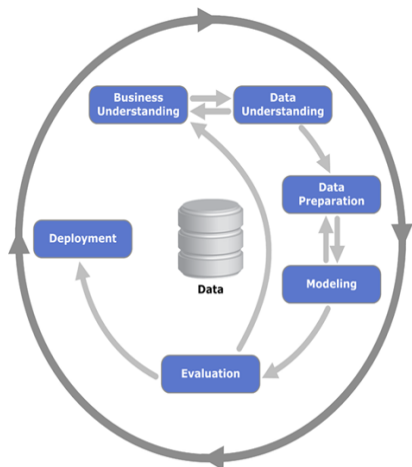
متدولوژی‌های متفاوتی برای داده کاوی از قبل سال ۱۹۹۶ ارائه شده است که در شکل‌های زیر مهم‌ترین آن‌ها آورده شده است:

^{۱۹} Data Mining

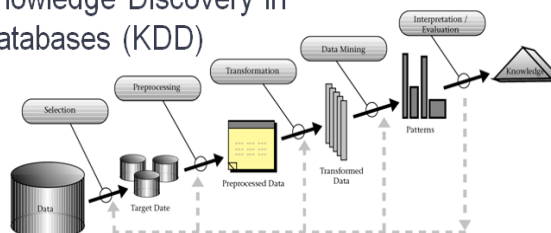
^{۲۰} Knowledge Discovery From Data | KDD

From Data to Insight

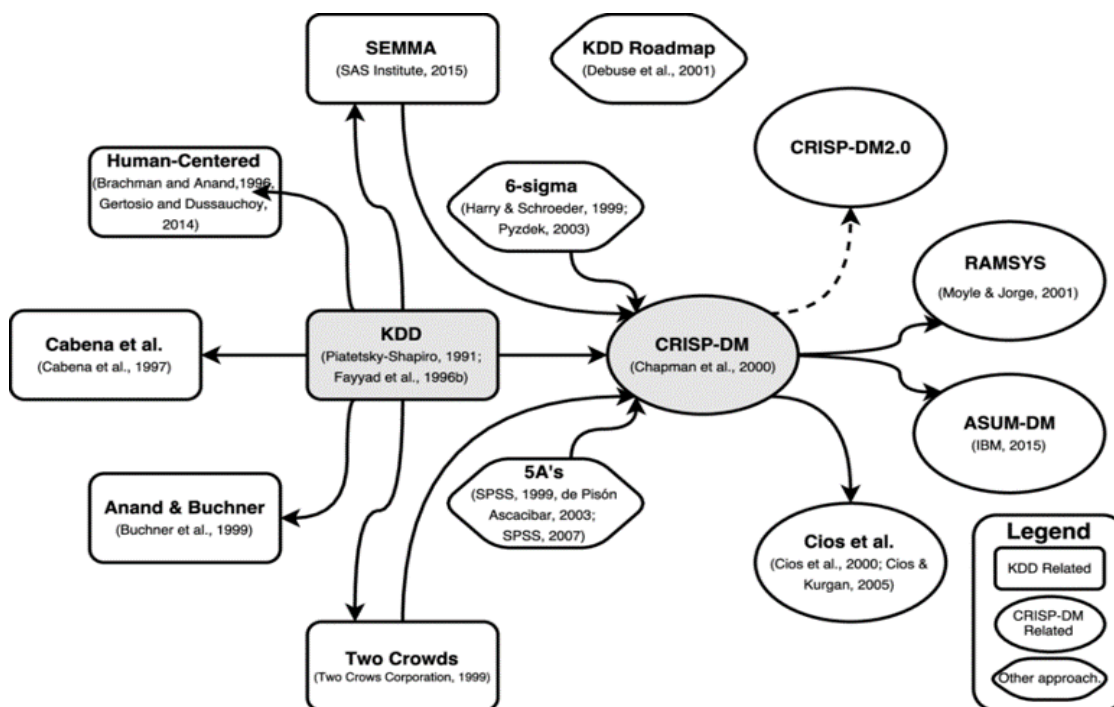
Cross Industry Standard Process
for Data Mining (CRISP-DM)
(IBM, Teradata, Daimler AG, NCR Corporation and OHRA)



Knowledge Discovery in
Databases (KDD)



شکل (۲-۴) متدولوژی‌های داده‌کاوی



شکل (۳-۴) متدولوژی‌های داده‌کاوی

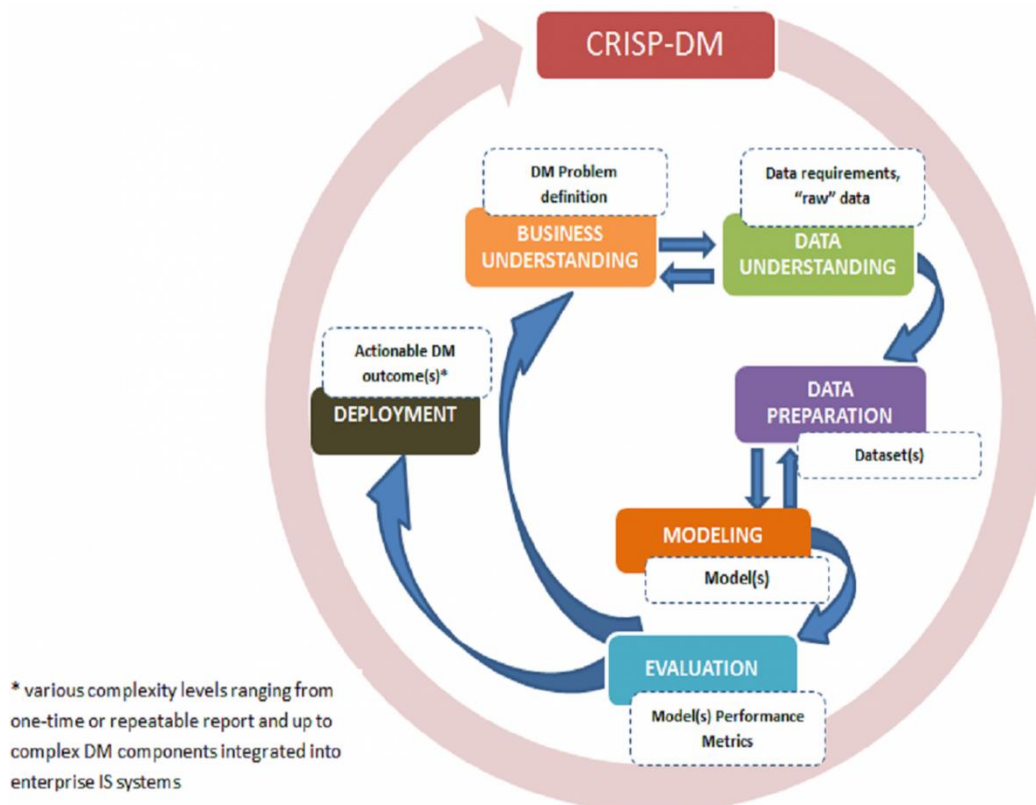
برخی از معروف‌ترین مدل‌های داده‌کاوی مدل‌های مبتنی بر CRISP-DM از شرکت IBM و مدل‌های مبتنی بر KDD است که در ادامه باتوجه به اهمیت مدل CRISP-DM توضیح مختصری در مورد آن داده خواهد شد.

۴-۵- مدل CRISP-DM

روش داده‌کاوی مبتنی بر CRISP-DM در ۶ گام زیر صورت می‌گیرد:

۱. درک کسب‌وکار
۲. بررسی و درک داده‌ها
۳. آماده‌سازی داده‌ها
۴. مدل‌سازی
۵. تست و ارزیابی مدل

توسعه مدل نهایی و استقرار



شکل (۴-۴) مدل داده‌کاوی CRISP-DM

۴-۵-۱- مرحله اول: درک کسب و کار

کاربران برای اتخاذ تصمیم‌های مناسب در هنگام ایجاد مدل‌های داده‌کاوی باید به درک صحیحی از داده‌ها برسند. در این مرحله مواردی همچون الزامات مربوط به کسب و کار، تعریف چارچوب مسئله، تعریف معیارهای مورد استفاده برای ارزیابی مدل و تعریف اهداف مشخص برای پروژه داده‌کاوی صورت می‌پذیرد.

۴-۵-۲- مرحله دوم: بررسی و درک داده‌ها

متخصص داده‌کاوی، داده‌های ثبت شده در کسب و کار کارفرما را از وی درخواست می‌کند و به بررسی داده‌ها می‌پردازد. متخصص داده‌کاوی باتوجه به حجم و کیفیت داده‌ها، مسئله طرح شده در مرحله قبل را تعدیل می‌کند تا نتیجه‌ی پروسه‌ی داده‌کاوی واقع‌بینانه‌تر به شبه‌طور به طور خلاصه می‌توان دیتاست‌ها را به موارد زیر تقسیم‌بندی نمود:

- کمی: همانند درجه حرارت و قد افراد
- کیفی: همانند دسته مدارک تحصیلی (دیپلم، لیسانس و...) یا گروه رنگ‌ها (زرد، قرمز و...)
- ترتیبی: چنین داده‌هایی دارای یک ترتیب طبیعی هستند؛ همانند مدارج تحصیلی (دبستان، راهنمایی، دبیرستان، کارشناسی، کارشناسی ارشد، دکتری)
- اسمی: اسمی دسته‌ها همانند وضعیت تأهل، جنسیت و رنگ‌ها
- عددی: داده‌های عددی خود به دودسته فاصله‌ای و نسبتی تقسیم می‌شوند. داده‌های فاصله‌ای بر اساس مقیاس واحدهایی با اندازه برابر اندازه‌گیری می‌شوند. مقادیر ویژگی‌های عددی دارای ترتیب هستند و می‌توانند مثبت، صفر و یا منفی باشند. یک داده نسبتی، خصیصه عددی دارای یک صفر مطلق است. اگر اندازه‌ها نسبتی باشند، می‌توان از نسبت مقادیر با یکدیگر سخن گفت. به‌علاوه، مقادیر قابل مرتب‌سازی شدن هستند و می‌توان تفاضل بین آن‌ها، میانگین و مد را محاسبه نمود.

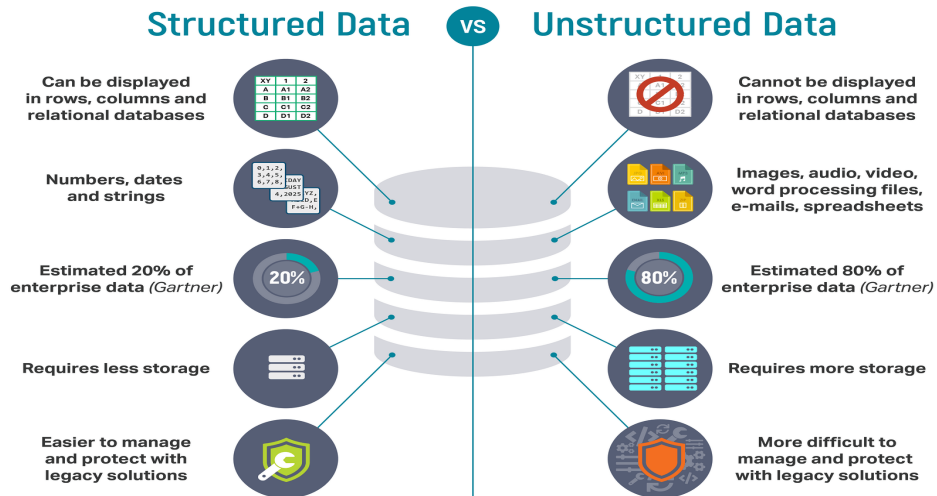
در بسیاری از مباحث داده‌کاوی، یادگیری ماشین و کلان‌داده‌ها، داده‌ها را می‌توان به دودسته تقسیم‌بندی کرد:

۱- داده‌های ساختاریافته^{۲۱}

۲- داده‌های غیرساختاریافته^{۲۲}

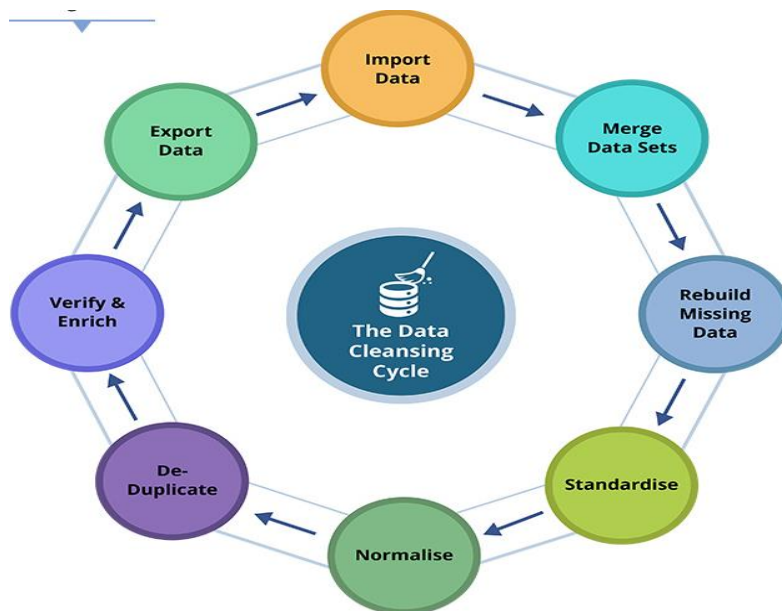
^{۲۱} Structured Data

^{۲۲} Un Structured Data

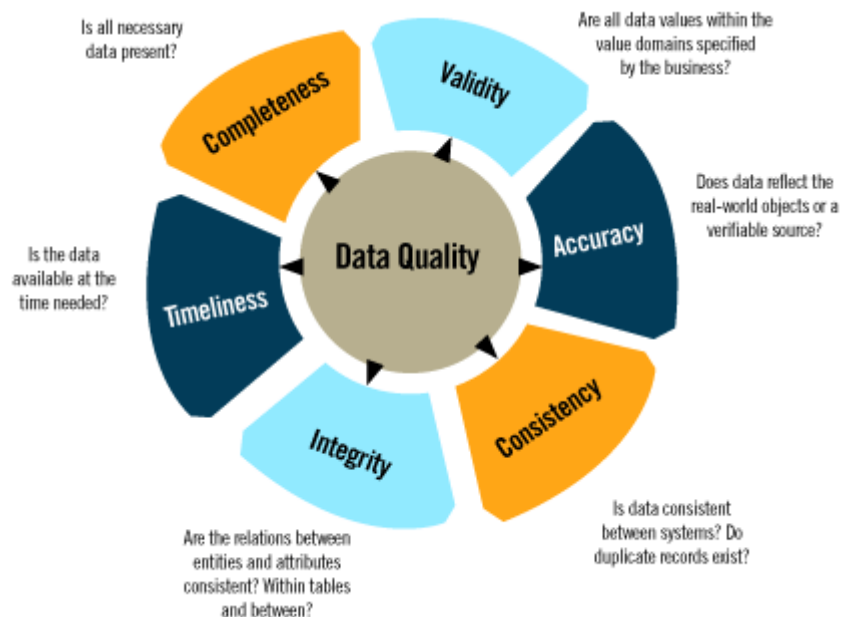


شکل (۴-۵) داده ساختاریافته و غیرساختاریافته

۴-۵-۳- مرحله سوم: آماده سازی یا پیش پردازش داده ها
 این امکان وجود دارد که داده ها در سراسر سازمان توزیع شده و در قالب های مختلف ذخیره گردند و یا اینکه ممکن است شامل تناقضات و ناسازگاری هایی از جمله ورودی های نادرست یا ازدست رفته باشند.



شکل (۴-۶) پاک سازی داده



شکل (۴-۷) اعتبارسنجی داده

۴-۵-۴- مرحله چهارم: مدل سازی

قدم چهارم مدل سازی داده های آماده سازی شده است. باتوجه به متدهای متفاوت، مدل های متفاوتی ساخته می شود و بهترین مدل ها از نظر متخصص داده کاوی انتخاب می شود.

۴-۵-۵- مرحله پنجم: تست و ارزیابی مدل

پیش از پیاده سازی مدل در محیط عملیاتی باید نحوه عملکرد آن مورد بررسی قرار گیرد. به علاوه در هنگام تهیه مدل معمولاً باید چندین مدل با پیکربندی های متفاوت ارائه شوند تا پس از تست نمودن آنها بتوان به مدلی دست یافت که بهترین نتیجه را در ارتباط با مشکلات و داده ها فراهم آورد. مدل های ساخته شده تست و ارزیابی می شوند و بهترین مدل از نظر مسئله طرح شده در مرحله یک، انتخاب می شود.

۴-۵-۶- مرحله ششم: توسعه مدل نهایی و استقرار

پس از استقرار Mining Model در یک محیط عملیاتی می توان عملکردهای بسیاری را باتوجه به نیازها اجرا نمود. در زیر به برخی از این عملکردها اشاره می شود. استفاده از مدل ها برای فرایندهای پیش بینی که ممکن است در مراحل بعدی برای اتخاذ تصمیمات در کسب و کار نیز به کار گرفته شود.

- انجام Query های محتوا به منظور بازایی اطلاعات آماری، قواعد یا فرمول های مربوط به مدل ها
- جای گذاری مستقیم عملکرد داده کاوی در برنامه های کاربردی

- ارائه گزارشی که امکان Query نمودن مستقیم در مدل داده‌کاوی موجود را برای کاربران فراهم می‌کند.

۴-۶- یادگیری ماشین^{۲۳}

یادگیری ماشین شاخه‌ای از هوش مصنوعی و علوم کامپیوتر است که بر استفاده از داده‌ها و الگوریتم‌ها برای تقلید از روشی که انسان‌ها یاد می‌گیرند، تمرکز دارد و به تدریج دقت آن را بهبود می‌بخشد. یادگیری ماشین جزء مهمی از حوزه روبه‌رشد علم داده است که از طریق استفاده از روش‌های آماری، الگوریتم‌ها، برای دسته‌بندی یا پیش‌بینی و کشف بینش‌های کلیدی در پروژه‌های داده‌کاوی آموزش داده می‌شوند. این بینش‌ها متعاقباً تصمیم‌گیری را در برنامه‌ها و کسب‌وکارها هدایت می‌کنند و به طور ایده‌آل بر معیارهای رشد کلیدی تأثیر می‌گذارند. الگوریتم‌های یادگیری ماشین از داده‌های ساختاریافته و برچسب‌گذاری شده برای پیش‌بینی استفاده می‌کنند. به این معنی که ویژگی‌های خاصی از داده‌های ورودی برای مدل تعریف شده و در جداول سازمان‌دهی می‌شوند. این لزوماً به این معنی نیست که از داده‌های بدون ساختار استفاده نمی‌کنند. این فقط به این معنی است که معمولاً برای سازماندهی داده‌ها در قالبی ساختاریافته، داده‌ها فرایند پیش‌پردازش را طی می‌کنند.

۴-۶-۱- انواع یادگیری ماشین

- یادگیری با نظارت^{۲۴}
- یادگیری بدون نظارت^{۲۵}
- یادگیری تقویتی^{۲۶}

۴-۶-۱-۱- الگوریتم‌های یادگیری ماشین با نظارت

این نوع از یادگیری، یک نوع یادگیری ماشین است که در آن الگوریتم از داده‌های برچسب‌دار یاد می‌گیرد. داده برچسب‌گذاری شده، به معنای مجموعه داده‌ای است که متغیر هدف مربوطه آن از قبل مشخص است. یادگیری با نظارت دو نوع دارد.

- **دسته‌بندی:** در این نوع از الگوریتم‌ها، کلاس مجموعه داده بر اساس متغیر ورودی مستقل پیش‌بینی می‌شود. کلاس مقادیر مقوله‌ای گسسته است. مثلاً تصویر حیوان گربه یا سگ است.
- **رگرسیون:** در این نوع از الگوریتم‌ها متغیرهای خروجی پیوسته بر اساس متغیر ورودی مستقل پیش‌بینی می‌شود. برای مثال پیش‌بینی قیمت مسکن بر اساس پارامترهای مختلف مانند سن خانه، فاصله از جاده اصلی، موقعیت مکانی، مساحت و غیره.

^{۲۳} Machine Learning

^{۲۴} Supervised Learning

^{۲۵} Unsupervised Learning

^{۲۶} Reinforcement Learning

۴-۶-۱-۲- الگوریتم‌های یادگیری ماشین بدون نظارت

در یادگیری بدون نظارت، الگوریتم باید خود به تنهایی به دنبال ساختارهای جالب موجود در داده‌ها باشد. به بیان ریاضی، یادگیری بدون نظارت مربوط به زمانی است که در مجموعه داده فقط متغیرهای ورودی وجود داشته باشند و هیچ متغیر داده خروجی موجود نباشد. به این نوع یادگیری، بدون نظارت گفته می‌شود. زیرا برخلاف یادگیری با نظارت، هیچ پاسخ صحیح داده شده‌ای وجود ندارد و ماشین خود باید به دنبال پاسخ باشد.

به بیان دیگر، هنگامی که الگوریتم برای کارکردن از مجموعه داده‌ای بهره گیرد که فاقد داده‌های برچسب دار (متغیرهای خروجی) است، از مکانیزم دیگری برای یادگیری و تصمیم‌گیری استفاده می‌کند. به چنین نوع یادگیری، بدون نظارت گفته می‌شود. یادگیری بدون نظارت قابل تقسیم به مسائل خوشه‌بندی و انجمنی است.

- **قوانین انجمنی:** یک مسئله یادگیری هنگامی قوانین انجمنی محسوب می‌شود که هدف کشف کردن قواعدی باشد که بخش بزرگی از داده‌ها را توصیف می‌کنند. مثلاً شخصی که کالای الف را خریداری کند، تمایل به خرید کالای ب نیز دارد.

- **خوشه‌بندی:** یک مسئله هنگامی خوشه‌بندی محسوب می‌شود که قصد کشف گروه‌های ذاتی (داده‌هایی که ذاتاً در یک گروه خاص می‌گنجند) در داده‌ها وجود داشته باشد. مثلاً، بخش‌بندی مشتریان بر اساس رفتار خرید آن‌ها.

۴-۶-۱-۳- یادگیری تقویتی

یک برنامه رایانه‌ای که با محیط پویا در تعامل است باید به هدف خاصی دست یابد (مانند بازی کردن با یک رقیب یا راندن خودرو). این برنامه بازخوردهایی را با عنوان پاداش‌ها و تنبیه‌ها فراهم و فضای مسئله خود را بر همین اساس هدایت می‌کند. با استفاده از یادگیری تقویتی، ماشین می‌آموزد که تصمیمات مشخصی را در محیطی که دائم در معرض آزمون و خطا است، اتخاذ کند.

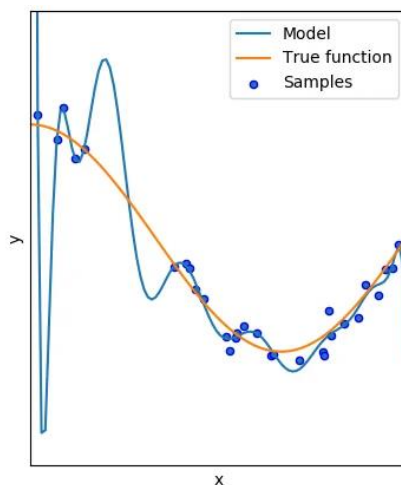
۴-۶-۲- بیش برآزش^{۲۷}، کم برآزش^{۲۸} و برآزش مناسب

مدل بیش برآزش، مدلی بسیار پیچیده برای داده‌ها است. به این معنی که در تحلیل رگرسیونی، مدلی با بیشترین پارامترها ایجاد می‌شود. در چنین حالتی، مدل با تغییرات جهشی سعی در پوشش داده‌های حاصل از نمونه و حتی مقدارهای نویز می‌کند. درحالی که چنین مدلی باید منعکس کننده رفتار جامعه باشد. در این گونه موارد، اگر مدل رگرسیون به دست آمده، برای پیش‌بینی نمونه دیگری به کار رود، مقدارهای پیش‌بینی شده اصلاً مناسب به نظر نخواهند رسید.

^{۲۷} Overfitting

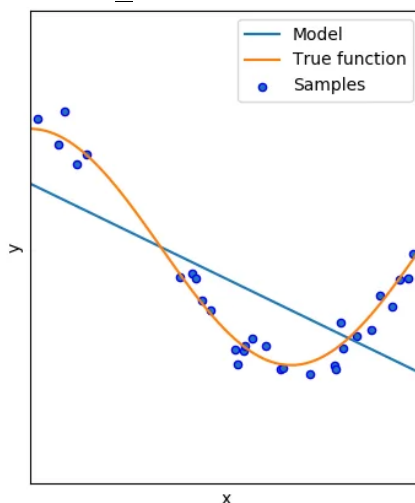
^{۲۸} Underfitting

در تصویر زیر، نمودار حاصل از بیش برازش روی داده‌های حاصل از نمونه دیده می‌شود. خط آبی، نشان‌دهنده منحنی برازش شده روی داده‌ها است و خط نارنجی تابعی است که مدل واقعی جامعه آماری را نشان می‌دهد. نقاط آبی‌رنگ نیز نمونه‌های تصادفی از جامعه آماری را نشان می‌دهند. در مدل بیش برازش، نقطه‌های حاصل از نمونه بهترین برازش را دارند و خط آبی تقریباً از همه آن‌ها عبور کرده است.



نمودار (۱-۴) منحنی بیش برازش بر اساس چندجمله‌ای مرتبه ۱۵

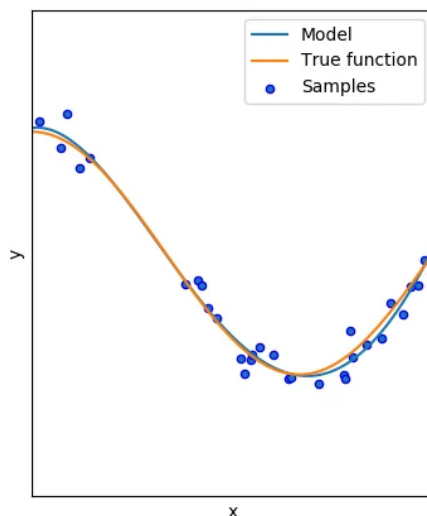
همچنین در زمانی که پارامترهای مدل رگرسیونی به صورت کم برازش برآورد می‌شوند، جانب احتیاط حفظ شده و مدل سعی می‌کند با کمترین پارامترها، عمل برازش را انجام دهد. در نتیجه خطای حاصل از این مدل حتی بر اساس نمونه‌های به کاررفته نیز بسیار زیاد است. در تصویر زیر، یک نمونه از مدل رگرسیونی کم برازش دیده می‌شود. درجه منحنی به کاررفته در این حالت ۱ است که معادله خط محسوب می‌شود.



نمودار (۲-۴) منحنی کم برازش بر اساس چندجمله‌ای مرتبه ۱

انتظار ما از یک تحلیل رگرسیون مناسب، ایجاد مدلی است که نه تنها بتواند برای داده‌های مربوط به نمونه، برازش مناسب را انجام دهد، بلکه برای داده‌هایی جدید نیز امکان برآورد مناسب وجود داشته باشد.

همان‌طور که در تصویر زیر دیده می‌شود، مدل مناسب دارای خطای کوچکی است و قابلیت پیش‌بینی برای داده‌های جدید را دارد.



نمودار (۳-۴) منحنی برازش مناسب بر اساس چندجمله‌ای مرتبه ۴

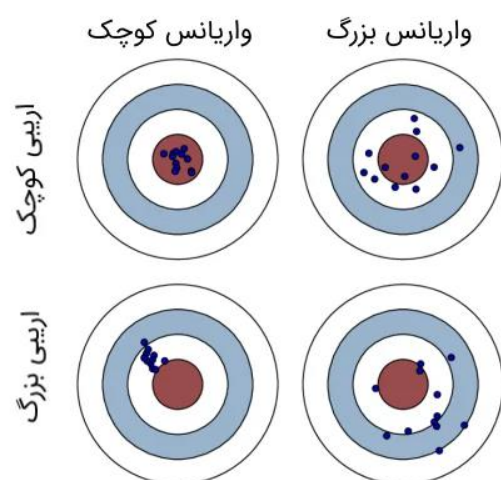
۳-۶-۴- موازنه واریانس و بایاس^{۲۹}

خطای بایاس: بایاس در واقع میزان اختلاف نقاط پیش‌بینی شده از متغیر هدف واقعی است. وجود فرضیه‌های مختلف روی مدل و الگوریتم یادگیری منجر به ایجاد خطای اریبی می‌شود. بزرگ‌بودن اریبی می‌تواند الگوریتم یا مدل آماری را از کشف روابط بین ویژگی‌ها و متغیر پاسخ باز دارد. اغلب بزرگ‌بودن خطای اریبی، منجر به کم‌برازش می‌شود.

خطای واریانس: واریانس میزان پراکندگی نقاط را نشان می‌دهد. هر چه واریانس بیشتر باشد، پراکندگی داده‌ها بیشتر است. حساسیت زیاد مدل با تغییرات کوچک روی داده‌های آموزشی، نشانگر وجود واریانس زیاد است. این امر نشانگر آن است که اگر مدل آموزش داده‌شده را روی داده‌های آزمایشی به کار گیریم، نتایج حاصل با داده‌های واقعی فاصله زیادی خواهند داشت. متأسفانه افزایش واریانس در این حالت منجر به مدل‌بندی مقادیر نویز^{۳۰} شده و به‌جای پیش‌بینی صحیح، دچار پیچیدگی و مشکل بیش‌برازش می‌شود. شکل زیر مصورسازی مفهوم موازنه واریانس و بایاس را نشان می‌دهد.

^{۲۹} Bias-Variance Tradeoff

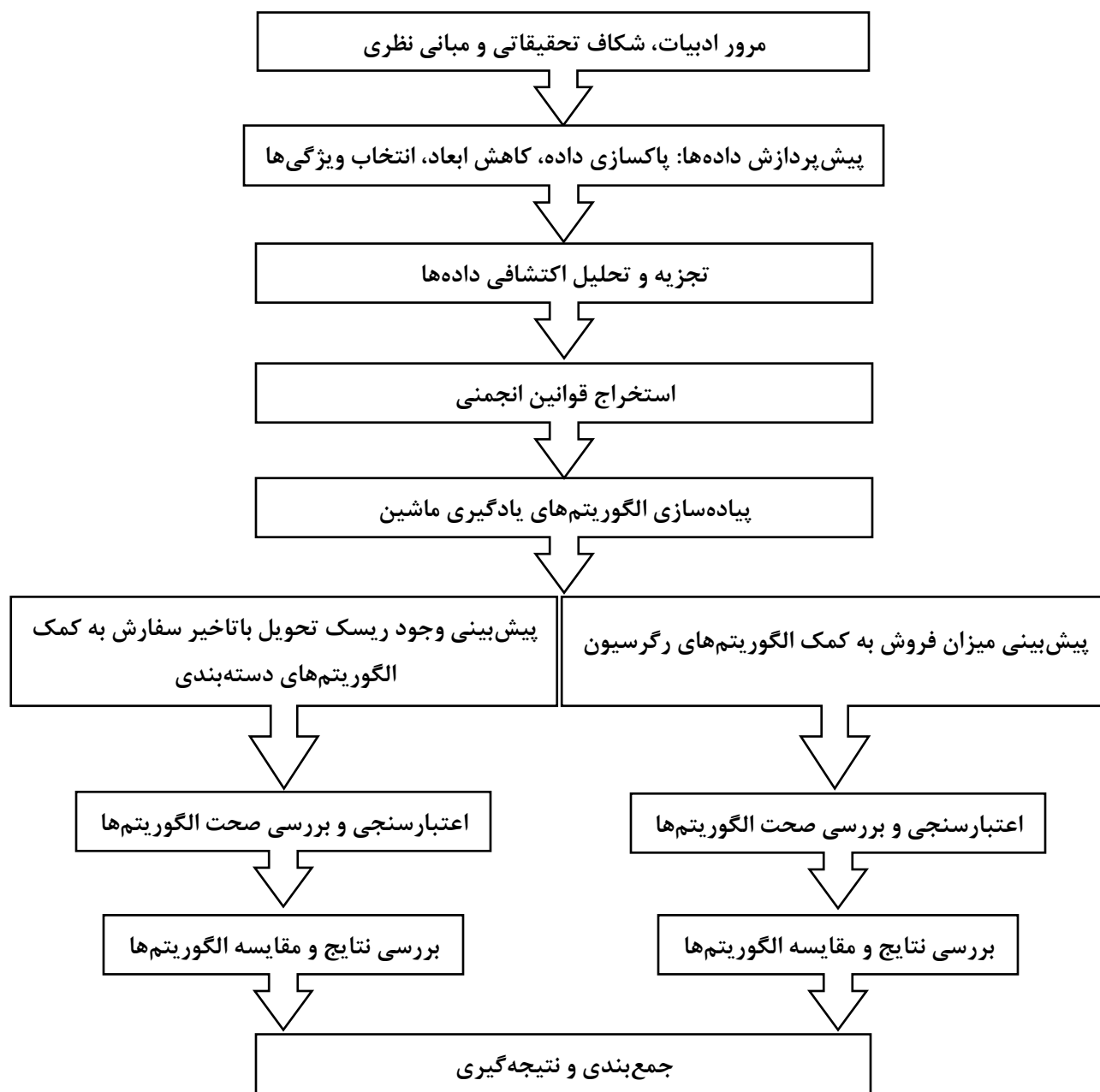
^{۳۰} Noise



شکل (۴-۸) مصورسازی موازنه واریانس و بایاس برای برآوردگر

۵- روش تحقیق

روش تحقیق این مطالعه، در شکل زیر آمده است.



شکل (۵-۱) روش تحقیق مطالعه

۵-۱- پیش‌پردازش داده‌ها: پاک‌سازی داده، کاهش ابعاد، انتخاب ویژگی‌ها

مجموعه داده توسط شرکت دیتاکو منتشر شده است. سه سطر از داده‌ها که غلط وارد شده‌اند، حذف گردید. از ستون‌های تاریخ و زمان سفارش و تاریخ و زمان حمل‌ونقل، متغیرهای زمان به تفکیک استخراج

گردید و به ستون‌های مجموعه داده اضافه شدند و فرایند استخراج ویژگی انجام شد. ستون‌هایی که تعداد زیادی از داده‌هایشان وجود نداشت، حذف شدند و در نهایت سطرهایی که داده‌های خالی داشتند، حذف گردید.

۵-۲- تجزیه و تحلیل اکتشافی داده‌ها

در این بخش، تجزیه تحلیل فراوانی روی ستون‌های (ویژگی‌های) موجود در مجموعه داده صورت گرفته است و با استخراج دانش، پیشنهاداتی برای بهبود ارائه گردیده است. همچنین همبستگی و ارتباط بین ستون‌ها نیز در این بخش مورد بررسی قرار گرفته است.

۵-۳- استخراج قوانین انجمنی^{۳۱}

استخراج قوانین انجمنی، تکنیکی است که برای کشف روابط پنهان بین متغیرها در مجموعه داده‌های بزرگ استفاده می‌شود. این یک روش محبوب در داده کاوی و یادگیری ماشین است و کاربردهای گسترده‌ای در زمینه‌های مختلف مانند تجزیه و تحلیل سبد بازار، تقسیم‌بندی مشتریان و کشف تقلب دارد. هدف از استخراج قوانین انجمنی، کشف قوانینی است که روابط بین متغیرهای مختلف در مجموعه داده را توصیف می‌کند.

برای مثال، مجموعه داده‌ای از معاملات در یک فروشگاه مواد غذایی را در نظر بگیرید. استخراج قوانین انجمنی می‌تواند برای شناسایی روابط بین اقلامی که اغلب با هم خریداری می‌شوند، استفاده شود. برای مثال، قانون «اگر مشتری نان بخرد، احتمالاً شیر هم می‌خرد.» یک قانون انجمنی است که می‌تواند از این مجموعه داده استخراج شود. ما می‌توانیم از چنین قوانینی برای اطلاع از تصمیم‌گیری‌ها در مورد چیدمان فروشگاه، قرارگیری محصول و بازاریابی استفاده کنیم.

الگوریتم‌های مختلفی برای استخراج قوانین انجمنی وجود دارند. در ادامه به پرکاربردترین آن‌ها اشاره می‌شود.

- **الگوریتم Apriori:** الگوریتم Apriori یکی از پرکاربردترین الگوریتم‌ها برای استخراج قوانین انجمنی است. این الگوریتم، ابتدا مجموعه موارد پرتکرار در مجموعه داده را شناسایی می‌کند (مجموعه‌هایی که در تعداد معینی از رکوردها ظاهر می‌شوند). سپس از این مجموعه موارد پرتکرار برای تولید قوانین انجمنی استفاده می‌کند. الگوریتم Apriori از یک رویکرد پایین‌به‌بالا استفاده می‌کند که از موارد جداگانه شروع می‌شود و به تدریج به مجموعه‌های موارد پیچیده‌تر می‌رسد.

^{۳۱} Association Rule Mining

- **الگوریتم FP-Growth**^{۳۲}: الگوریتم رشد الگوی پرتکرار، یکی دیگر از الگوریتم‌های محبوب

برای استخراج قوانین انجمنی است که با ساختن یک ساختار درخت‌مانند به نام FP-tree کار می‌کند که مجموعه موارد پرتکرار در مجموعه داده را رمزگذاری می‌کند. سپس از FP-tree برای ایجاد قوانین انجمنی به روشی مشابه الگوریتم Apriori استفاده می‌شود. الگوریتم رشد الگوی پرتکرار به‌طور کلی سریع‌تر از الگوریتم Apriori است.

- **الگوریتم ECLAT**^{۳۳}: الگوریتم خوشه‌بندی کلاس هم ارز و پیمایش شبکه از بالا به پایین،

نوعی از الگوریتم Apriori است که از رویکرد بالا به پایین به جای رویکرد از پایین به بالا استفاده می‌کند. با تقسیم موارد به کلاس‌های معادل بر اساس پشتیبانی آن‌ها (تعداد رکوردهایی که در آن‌ها ظاهر می‌شوند) کار می‌کند. سپس قوانین انجمنی با ترکیب این کلاس‌های هم ارزی در یک ساختار شبکه مانند ایجاد می‌شود. این یک نسخه کارآمدتر و مقیاس‌پذیرتر از الگوریتم Apriori است.

در این مطالعه ما از الگوریتم Apriori برای استخراج قوانین انجمنی استفاده کرده‌ایم. در ادامه نحوه عملکرد این الگوریتم را بیان می‌کنیم.

۵-۳-۱- الگوریتم Apriori

این الگوریتم، با تنظیم حداقل آستانه پشتیبانی^{۳۴} شروع می‌شود. این عدد حداقل تعداد دفعاتی است که یک مورد باید در پایگاه داده رخ دهد تا بتوان آن را به عنوان مجموعه موارد پرتکرار در نظر گرفت. سپس الگوریتم هر مجموعه مواردی را که حداقل آستانه پشتیبانی را برآورده نمی‌کنند، فیلتر می‌کند.

سپس الگوریتم لیستی از تمام ترکیبات ممکن از مجموعه موارد پرتکرار ایجاد می‌کند و تعداد دفعاتی که هر ترکیب در پایگاه داده ظاهر می‌شود را می‌شمارد. در ادامه الگوریتم فهرستی از قوانین مرتبط را بر اساس ترکیبات پرتکرار مجموعه موارد تولید می‌کند.

قدرت^{۳۵} قانون انجمنی با استفاده از معیار اطمینان^{۳۶} اندازه‌گیری می‌شود که احتمال وجود مورد ب با توجه به وجود مورد الف است. سپس الگوریتم قوانین انجمنی را که حداقل آستانه اطمینان را برآورده نمی‌کند، فیلتر می‌کند. از این قوانین به عنوان قوانین انجمنی قوی یاد می‌شود. در نهایت، الگوریتم لیستی از قوانین مرتبط قوی را به عنوان خروجی برمی‌گرداند. در ادامه به معیارهای ارزیابی و تحلیل قوانین انجمنی می‌پردازیم.

^{۳۲} Frequent Pattern Growth

^{۳۳} Equivalence Class Clustering and bottom-up Lattice Traversal

^{۳۴} Minimum Support Threshold

^{۳۵} Strength

^{۳۶} Confidence

۵-۳-۲- معیارهای ارزیابی قوانین انجمن

در استخراج قوانین انجمنی، معمولاً از چندین معیار برای ارزیابی کیفیت و اهمیت قوانین کشف شده استفاده می‌شود. این معیارها را می‌توان برای ارزیابی کیفیت و اهمیت قوانین مرتبط و انتخاب مناسب‌ترین قوانین برای یک کاربرد خاص مورد استفاده قرارداد.

تفسیر نتایج معیارهای استخراج قواعد انجمنی مستلزم درک معنا و مفاهیم هر معیار و همچنین نحوه استفاده از آن‌ها برای ارزیابی کیفیت و اهمیت قوانین مرتبط کشف شده است. در اینجا چند دستورالعمل برای تفسیر نتایج معیارهای استخراج قانون انجمنی اصلی آورده شده است.

- **Support:** پشتیبانی معیاری است که نشان می‌دهد یک مورد یا مجموعه موارد به دفعات در مجموعه داده ظاهر می‌شود و به شکل زیر محاسبه می‌گردد. پشتیبانی زیاد نشان می‌دهد که یک مورد یا مجموعه موارد در مجموعه داده مشترک است، در حالی که پشتیبانی کم نشان‌دهنده نادر بودن آن است.

$$Support(\{X\} \rightarrow \{Y\}) = \frac{\text{Transactions containing both } X \text{ and } Y}{\text{Total number of transactions}}$$

- **Confidence:** معیاری برای سنجش قدرت ارتباط بین دو مورد است. به عنوان تعداد رکوردهای حاوی هر دو مورد تقسیم بر تعداد رکوردهای حاوی اولین مورد محاسبه می‌شود. اطمینان بالا نشان می‌دهد که وجود مورد اول یک پیش‌بینی‌کننده قوی برای حضور مورد دوم است.

$$Confidence(\{X\} \rightarrow \{Y\}) = \frac{\text{Transactions containing both } X \text{ and } Y}{\text{Transactions containing } X}$$

- **Lift:** این معیار اندازه‌گیری قدرت ارتباط بین دو مورد با در نظر گرفتن فراوانی هر دو مورد در مجموعه داده است. به عنوان معیار اطمینان تقسیم بر معیار پشتیبانی مورد دوم محاسبه می‌شود. این معیار، برای مقایسه قدرت ارتباط بین دو مورد با قدرت مورد انتظار انجمن در صورتی که موارد مستقل باشند، استفاده می‌شود. مقدار بیشتر از ۱ نشان می‌دهد که ارتباط بین دو مورد قوی‌تر از حد انتظار بر اساس فراوانی اقلام است. این نشان می‌دهد که این ارتباط ممکن است معنی‌دار باشد و ارزش بررسی بیشتر را داشته باشد. مقدار کمتر از ۱ نشان می‌دهد که ارتباط ضعیف‌تر از حد انتظار است و احتمالاً کمتر قابل توجه باشد.

$$Lift(\{X\} \rightarrow \{Y\}) = \frac{(\text{Transactions containing both } X \text{ and } Y) / (\text{Transactions containing } X)}{\text{Fraction of transactions containing } Y}$$

۵-۴- پیاده‌سازی الگوریتم‌های یادگیری ماشین

۵-۴-۱- آماده‌سازی مجموعه داده برای پیاده‌سازی الگوریتم‌های یادگیری ماشین

از آنجایی که ستون‌هایی که واریانس کمی دارند در مدل‌های پیش‌بینی تأثیر چندانی ندارند، برای کاهش زمان پردازش مدل، از مجموعه داده حذف شدند. همچنین ستون‌های تاریخ و زمان سفارش و تاریخ و زمان

حمل و نقل به علت اینکه نوع داده‌های آن‌ها از نوع عددی^{۳۷} نمی‌باشند نیز از مجموعه داده حذف گردیدند. زیرا الگوریتم‌های یادگیری ماشین تنها روی داده‌های عددی قابل اجرا هستند.

برای تبدیل سایر ستون‌های غیر عددی طبقه‌بندی شده^{۳۸}، به ستون‌های عددی از دو روش کدگذاری طبقه‌ای^{۳۹} استفاده شده است. کدگذاری طبقه‌ای، فرایندی است که در آن داده‌های طبقه‌بندی به داده‌های عددی تبدیل می‌شوند. تکنیک‌های کدگذاری طبقه‌بندی زیادی وجود دارد، در این مطالعه، ما از دو روش کدگذاری وان‌هات^{۴۰} و کدگذاری هَش^{۴۱} استفاده کرده‌ایم.

کدگذاری وان‌هات، تکنیکی است که تمام عناصر یک ستون طبقه‌بندی را به ستون‌های جدیدی تبدیل می‌کند که با ۰ یا ۱ نشان داده می‌شوند تا وجود مقدار دسته را نشان دهد. جدول زیر مثالی از فرایند این روش را نشان می‌دهد.

جدول (۵-۱) فرایند کدگذاری وان‌هات

Original categorical column	One-Hot encoded columns		
Origin	Origin_USA	Origin_Japan	Origin_Europe
USA	1	0	0
Japan	0	1	0
Europe	0	0	1
USA	1	0	0
Europe	0	0	1

کدگذاری هَش، مبتنی بر تابع هَش انجام می‌شود و داده‌های طبقه‌بندی را به عددی تبدیل می‌کند. مزیت اصلی استفاده از کدگذاری هَش این است که می‌توان تعداد ستون‌های عدد مطلوب را کنترل کرد. این کار باعث می‌شود تا هنگام پردازش مجموعه داده توسط مدل‌های یادگیری ماشین، میزان حافظه^{۴۲} کمتری اشغال گردد و مدل در زمان کمتری اجرا شود.

در این مطالعه، ما فرایند کدگذاری وان‌هات را روی ستون‌های نوع معامله انجام شده، وضعیت تحویل سفارش‌ها، کشور محل سفارش، بازار محل تحویل سفارش و حالت حمل و نقل پیاده‌سازی کرده‌ایم و این

^{۳۷} Numerical

^{۳۸} Categorical

^{۳۹} Categorical Encoding

^{۴۰} One-Hot Encoding (OHE)

^{۴۱} Hash Encoding

^{۴۲} Memory

ستون‌های غیر عددی را به مقادیر عددی تبدیل کرده‌ایم. همچنین کدگذاری هش را روی ستون‌های شهر محل سفارش، بخش‌بندی مشتریان، ایالت محل سفارش، شهر مقصد سفارش، کشور مقصد سفارش، منطقه مقصد سفارش، ایالت مقصد سفارش و وضعیت سفارش پیاده‌سازی کرده‌ایم و تمام ستون‌های غیر عددی را به ستون‌های عددی تبدیل کرده‌ایم.

۵-۴-۲- تقسیم مجموعه داده به سه قسمت آموزشی^{۴۳}، اعتبارسنجی^{۴۴} و تست^{۴۵}

پس از آماده‌سازی مجموعه داده برای پیاده‌سازی الگوریتم‌های یادگیری ماشین، مجموعه داده را به سه قسمت آموزشی، اعتبارسنجی و تست تقسیم کرده‌ایم. ۷۰ درصد مجموعه داده، به عنوان مجموعه داده آموزشی، ۱۰ درصد برای اعتبارسنجی اولیه هر مدل و ۲۰ درصد از مجموعه داده به عنوان مجموعه داده تست در نظر گرفته شده‌اند که این مجموعه داده برای اندازه‌گیری صحت و دقت مدل‌ها کنار گذاشته می‌شود و اصطلاحاً توسط مدل‌ها دیده نمی‌شود و مدل در نهایت روی آن تست می‌گردد.

۵-۴-۳- پیاده‌سازی الگوریتم‌های یادگیری ماشین برای پیش‌بینی میزان فروش هر سفارش

از آنجایی که متغیر هدف ما متغیر پیوسته می‌باشد، برای پیش‌بینی میزان فروش هر سفارش، از الگوریتم‌های رگرسیون با نظارت زیر استفاده شده است.

۵-۴-۱- رگرسیون خطی^{۴۶}

رگرسیون خطی، نوعی از الگوریتم یادگیری ماشین با نظارت است که رابطه خطی بین یک متغیر وابسته و یک یا چند ویژگی مستقل را محاسبه می‌کند. هنگامی که تعداد ویژگی مستقل، ۱ باشد، به آن رگرسیون خطی تک‌متغیره و در مورد بیش از یک ویژگی، به عنوان رگرسیون خطی چندمتغیره شناخته می‌شود. هدف این الگوریتم، یافتن بهترین معادله خطی است که بتواند مقدار متغیر وابسته را بر اساس متغیرهای مستقل پیش‌بینی کند. این الگوریتم، معادله یک خط مستقیم را ارائه می‌دهد که نشان‌دهنده رابطه بین متغیرهای وابسته و مستقل است. در معادله ارائه شده، شیب خط نشان می‌دهد که متغیر وابسته برای یک واحد در متغیر(های) مستقل چقدر تغییر می‌کند.

رگرسیون خطی در بسیاری از زمینه‌های مختلف از جمله مالی، اقتصاد و روان‌شناسی برای درک و پیش‌بینی رفتار یک متغیر خاص استفاده می‌شود. به عنوان مثال، در امور مالی، رگرسیون خطی ممکن است برای درک رابطه بین قیمت سهام شرکت و درآمد آن یا برای پیش‌بینی ارزش آتی یک ارز بر اساس عملکرد گذشته آن استفاده شود.

^{۴۳} Train Dataset

^{۴۴} Validation Dataset

^{۴۵} Test Dataset

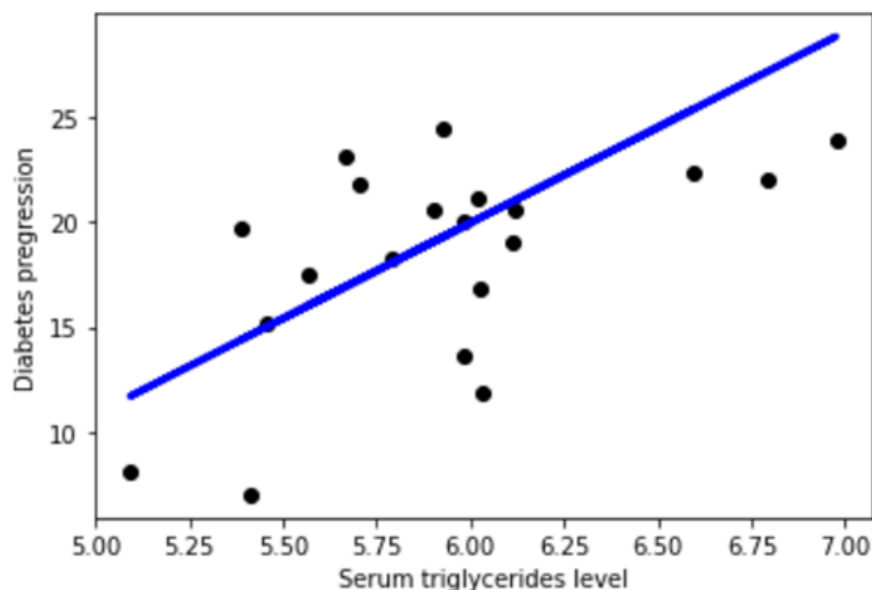
^{۴۶} Linear Regression

مدل رگرسیون خطی، معادله خطی را پیدا می‌کند که به بهترین شکل همبستگی متغیرهای مستقل را با متغیر وابسته توصیف می‌کند. این کار با برازش یک خط به داده‌ها با استفاده از روش حداقل مربعات به دست می‌آید. این خط سعی می‌کند مجموع مجذورهای باقی‌مانده را به حداقل برساند. باقی‌مانده فاصله بین خط و مقدار واقعی متغیر توضیحی است. در زیر نمونه‌ای از معادله رگرسیون خطی حاصل آمده است:

$$y = b_0 + b_1x_1 + b_2x_2 + \dots + e$$

در فرمول بالا، y متغیر وابسته و x_1 و x_2 و غیره متغیرهای مستقل هستند. ضرایب b_1 و b_2 و غیره همبستگی متغیرهای مستقل را با متغیر وابسته نشان می‌دهند. علامت ضرایب (+/-) مشخص می‌کند که آیا متغیر همبستگی مثبت یا منفی دارد. b_0 عرض از مبدایی است که مقدار متغیر وابسته را با فرض $x = 0$ بودن همه متغیرهای مستقل نشان می‌دهد.

در تصویر زیر، نمودار یک مدل رگرسیون خطی نشان داده شده است. این مدل رابطه بین متغیر وابسته، پیشروی دیابت، و متغیر مستقل، سطح تری گلیسیرید سرم^{۴۷} را توصیف می‌کند. همان‌طور که مشخص است یک همبستگی مثبت نشان داده شده است. این مثال یک مدل رگرسیون خطی با دو متغیر را نشان می‌دهد. اگرچه نمی‌توان مدل‌هایی با بیش از سه متغیر را تجسم کرد، اما عملاً یک مدل می‌تواند هر تعداد متغیر داشته باشد.



شکل (۵-۲) مثالی از پیاده‌سازی روش رگرسیون خطی

^{۴۷} Serum Triglycerides Level

۵-۴-۳-۲- رگرسیون ستیغی^{۴۸}

در مباحث مربوط به رگرسیون چندگانه، تعیین تعداد متغیرهای مستقلی که باید در مدل به کار گرفته شوند، یک مشکل محسوب می‌شود. با افزایش تعداد متغیرها، بیش برآزش^{۴۹} رخ داده و با کاهش آن‌ها نیز ممکن است با مسئله کم برآزش^{۵۰} مواجه شویم. در صورتی که مدل رگرسیونی دچار بیش برآزش شود، خطای آن برای برآورد مقدارهای جدید متغیر وابسته زیاد خواهد بود. در حالی که وجود متغیرهای کمتر از حد لازم در مدل (کم برآزش) واریانس مدل را افزایش می‌دهد؛ بنابراین با افزایش تعداد متغیرها مشکل هم‌خطی و بیش برآزش ظاهر شده و با کاهش آن‌ها، واریانس مدل افزایش خواهد یافت. یکی از روش‌های غلبه بر این مسائل در رگرسیون چندگانه، استفاده از مدل رگرسیون ستیغی است. از آنجایی که در زمانی که متغیرهای مدل، زیاد و یا هم‌خطی چندگانه وجود داشته باشد، واریانس برآوردگرها متورم شده و به شکل قله (ستیغ) در می‌آید، از همین روی، این روش رگرسیونی که بر این مشکل غلبه می‌کند، رگرسیون ستیغی نام‌گذاری شده است.

رگرسیون ستیغی، روشی برای تخمین ضرایب مدل‌های رگرسیون چندگانه در سناریوهایی است که متغیرهای مستقل همبستگی بالایی دارند. این الگوریتم در بسیاری از زمینه‌ها از جمله اقتصادسنجی، شیمی و مهندسی استفاده شده است. این روش به‌ویژه برای کاهش مشکل چندخطی بودن در رگرسیون خطی مفید است که معمولاً در مدل‌هایی با تعداد پارامترهای زیاد رخ می‌دهد. به‌طور کلی، این روش کارایی بهبودیافته‌ای را در مسائل تخمین پارامتر در ازای مقدار قابل تحملی سوگیری فراهم می‌کند.

رگرسیون ستیغی به‌عنوان یک راه‌حل ممکن برای عدم دقت برآوردگرهای حداقل مربع توسعه داده شد، زمانی که مدل‌های رگرسیون خطی دارای متغیرهای مستقل چندخطی (بسیار همبسته) هستند، با ایجاد یک برآوردگر رگرسیون ستیغی، یک تخمین دقیق‌تری از پارامترهای برجستگی ارائه می‌دهد، زیرا واریانس آن و برآوردگر میانگین مربع اغلب کوچک‌تر از برآوردگرهای حداقل مربعی است که قبلاً ارائه شده بود.

این الگوریتم، با جریمه کردن بزرگی ضرایب ویژگی‌ها و به حداقل رساندن خطا بین مشاهدات پیش‌بینی شده و واقعی کار می‌کند که به آن منظم‌سازی^{۵۱} می‌گویند. رگرسیون ستیغی، روش منظم‌سازی L_2 را انجام می‌دهد، یعنی جریمه‌ای معادل مجذور بزرگی ضرایب اضافه می‌کند. هدف این الگوریتم کمینه‌سازی هدف حداقل مربعات به‌اضافه حاصل ضرب ضریب ثابت α در مجموع مجذور ضرایب است که رابطه آن به‌صورت زیر است.

$$\text{Minimization objective} = \text{least squares objective} + \alpha * (\text{sum of square of coefficients})$$

^{۴۸} Ridge Regression

^{۴۹} Overfitting

^{۵۰} Underfitting

^{۵۱} Regularization

$$RIDGE : \quad RSS + \lambda \sum_{i=1}^n \beta_i^2$$

$$LASSO : \quad RSS + \lambda \sum_{i=1}^n |\beta_i|$$

که در آن مجموع مربعات باقی مانده (انحرافات پیش‌بینی شده از روی مقادیر تجربی واقعی) از رابطه زیر به دست می‌آید.

$$RSS = \sum_{i=1}^n (y_i - f(x_i))^2$$

۵-۴-۳-۴-۵- رگرسیون درخت تصمیم^{۵۸}

روش درخت تصمیم برای دسته‌بندی و پیش‌بینی مقادیر گسسته بکار می‌رود که در بخش‌های بعدی به تفصیل به آن پرداخته خواهد شد. به طور خلاصه، فرایند پیاده‌سازی الگوریتم درخت تصمیم از گره ریشه شروع می‌شود و توسط یک درخت منشعب دنبال می‌شود که در نهایت به یک گره برگ منتهی می‌شود که حاوی پیش‌بینی یا نتیجه نهایی الگوریتم است. ساخت درخت تصمیم معمولاً از بالا به پایین با انتخاب متغیری در هر مرحله که به بهترین نحو مجموعه موارد را تقسیم می‌کند، انجام می‌شود. هر زیر درخت از مدل درخت تصمیم را می‌توان به عنوان یک درخت دودویی نشان داد که در آن یک گره تصمیم بر اساس شرایط به دو گره تقسیم می‌شود. در درخت تصمیم متغیر هدف یا گره پایانی می‌تواند مقادیر پیوسته (معمولاً اعداد واقعی) باشد که به آن، درخت تصمیم رگرسیون گفته می‌شود.

۵-۴-۳-۵- رگرسیون جنگل تصادفی^{۵۹}

جنگل تصادفی یک روش یادگیری ترکیبی برای دسته‌بندی و رگرسیون می‌باشد که بر اساس ساختاری متشکل از شمار بسیاری درخت تصمیم، بر روی زمان آموزش و خروجی کلاس‌ها (دسته‌بندی) یا برای پیش‌بینی‌های هر درخت به شکل مجزا، کار می‌کنند. جنگل‌های تصادفی برای درختان تصمیم که در مجموعه آموزشی دچار بیش‌برازش می‌شوند، مناسب هستند. عملکرد جنگل تصادفی معمولاً بهتر از درخت تصمیم است، اما این بهبود عملکرد تا حدی به نوع داده هم‌بستگی دارد. هنگامی که متغیر هدف پیوسته است از رگرسیون جنگل تصادفی استفاده می‌شود.

^{۵۸} Decision Tree Regression

^{۵۹} Random Forest Regression

۵-۴-۳-۶- رگرسیون بیزی^{۶۰}

در رگرسیون خطی بیزی، میانگین یک پارامتر با مجموع وزنی سایر متغیرها مشخص می‌شود. هدف این نوع مدل‌سازی شرطی تعیین توزیع قبلی رگرسیون‌ها و همچنین سایر متغیرهایی است که تخصیص رگرسیون را توصیف می‌کنند و در نهایت امکان پیش‌بینی خارج از نمونه رگرسیون و مشروط به مشاهدات ضرایب رگرسیون را می‌دهد.

هنگامی که مجموعه داده دارای داده‌های بسیار کم یا ضعیف است، رگرسیون بیزی ممکن است بسیار مفید باشد. برخلاف روش‌های رگرسیون مرسوم که در آن خروجی تنها از یک عدد از هر ویژگی مشتق می‌شود، خروجی مدل رگرسیون بیزی از توزیع احتمال مشتق می‌شود.

۵-۴-۴- پیاده‌سازی الگوریتم‌های یادگیری ماشین برای پیش‌بینی وجود ریسک ارسال با تأخیر

برای پیش‌بینی وجود ریسک ارسال با تأخیر سفارش از الگوریتم‌های دسته‌بندی با نظارت زیر استفاده شده است. عدد ۱ به معنی وجود ریسک ارسال با تأخیر و عدد ۰ به معنی عدم وجود ریسک ارسال با تأخیر هستند.

۵-۴-۴-۱- درخت تصمیم‌گیری^{۶۱}

درخت تصمیم یک مدل سلسله‌مراتبی پشتیبانی تصمیم است که از یک مدل درخت‌مانند از تصمیمات و پیامدهای احتمالی آن‌ها، از جمله نتایج رویدادهای شانس، هزینه‌های منابع و مطلوبیت استفاده می‌کند. درخت‌های تصمیم معمولاً در تحقیقات عملیاتی، به‌ویژه در تجزیه و تحلیل تصمیم‌گیری برای کمک به شناسایی استراتژی که به احتمال زیاد به یک هدف می‌رسد، استفاده می‌شوند.

درخت تصمیم دارای اجزای زیر است:

- **گره اصلی^{۶۲}:** ویژگی کلیدی در مجموعه داده
- **گره داخلی^{۶۳}:** گره‌هایی که یک یال ورودی و دو یا چند یال خروجی دارند.
- **گره برگ^{۶۴}:** گره پایانی بدون یال خروجی

درخت تصمیم از یک گره اصلی شروع می‌شود و با بررسی شرایط مختلف و اختصاص آن به سایر گره‌ها ادامه می‌یابد. درخت تصمیم زمانی کامل می‌شود که تمام شرایط به یک گره برگ منتهی شوند. گره برگ حاوی برچسب طبقه‌بندی می‌باشد.

^{۶۰} Bayesian Regression

^{۶۱} Decision Tree

^{۶۲} Root Node

^{۶۳} Internal Node

^{۶۴} Leaf Node

برای تقسیم بهینه ویژگی‌ها دو روش وجود دارد:

- **روش شاخص جینی^{۶۵}:** ناخالصی جینی تعداد برچسب‌گذاری اشتباه هر عنصر مجموعه داده را هنگامی که به طور تصادفی برچسب‌گذاری می‌شود، اندازه‌گیری می‌کند. در شکل زیر فرمول شاخص جینی مشاهده می‌شود که در آن p_j احتمال کلاس j است. حداقل مقدار شاخص جینی $\frac{1}{2}$ است که زمانی اتفاق می‌افتد که گره خالص باشد، به این معنی که تمام عناصر موجود در گره از یک کلاس منحصر به فرد هستند؛ بنابراین، این گره دوباره تقسیم نخواهد شد. تقسیم بهینه توسط ویژگی‌هایی با شاخص جینی کمتر انتخاب می‌شود. علاوه بر این، زمانی که احتمال دو کلاس یکسان باشد، حداکثر مقدار (۰.۵) را دریافت می‌کند.

$$GiniIndex = 1 - \sum_j p_j^2$$

- **روش آنتروپی^{۶۶}:** آنتروپی معیاری از اطلاعات است که نشان‌دهنده بی‌نظمی ویژگی‌ها با متغیر هدف است. مشابه شاخص جینی، تقسیم بهینه توسط ویژگی با آنتروپی کمتر انتخاب می‌شود. مقدار آنتروپی زمانی حداکثر مقدار خود (۱) را به دست می‌آورد که احتمال دو کلاس یکسان باشد و هنگامی که یک گره خالص باشد، مقدار آنتروپی حداقل مقدار خود یعنی $\frac{1}{2}$ است. فرمول محاسبه آنتروپی به شکل زیر است که در آن p_j احتمال کلاس j است.

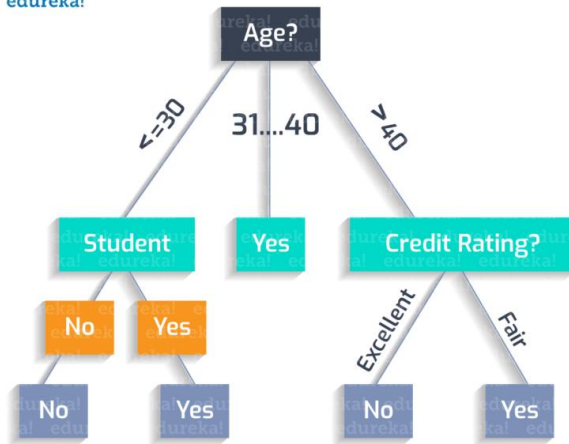
$$Entropy = - \sum_j p_j \cdot \log_2 \cdot p_j$$

شکل زیر نمونه‌ای از یک درخت تصمیم با تقسیم ویژگی‌ها با روش شاخص جینی را نشان می‌دهد که هدف آن پیش‌بینی خرید لپ‌تاپ توسط کاربر می‌باشد. همان‌طور که مشاهده می‌شود، ویژگی سن^{۶۷} به عنوان گره اصلی انتخاب شده است و سایر ویژگی‌ها در گره‌های داخلی قرار دارند و با بررسی شرایط مختلف، گره‌های برگ مشخص شده‌اند.

^{۶۵} Gini Index

^{۶۶} Entropy

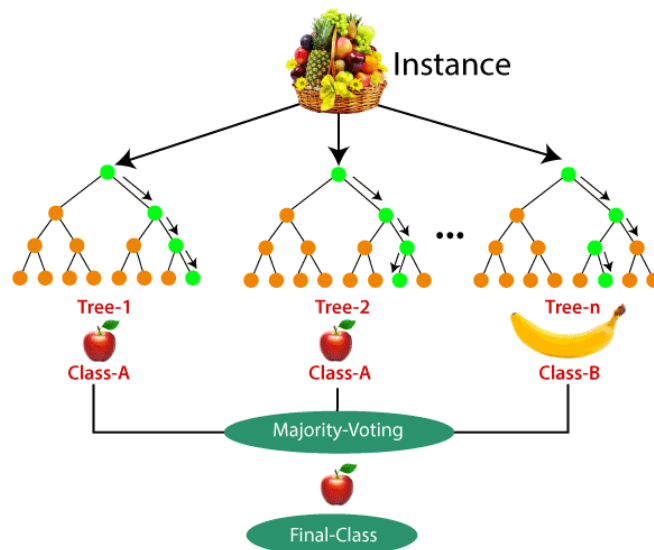
^{۶۷} Age



شکل (۳-۵) مثالی از پیاده‌سازی الگوریتم درخت تصمیم برای طبقه‌بندی داده‌ها

۵-۴-۲- جنگل تصادفی^{۶۸}

جنگل تصادفی یک روش یادگیری ترکیبی برای دسته‌بندی و رگرسیون می‌باشد که بر اساس ساختاری متشکل از شمار بسیاری درخت تصمیم، در زمان آموزش عمل می‌کند. عملکرد الگوریتم جنگل تصادفی معمولاً بهتر از الگوریتم درخت تصمیم است، اما این بهبود عملکرد تا حدی به نوع داده هم‌بستگی دارد. برای کاربرد دسته‌بندی، خروجی جنگل تصادفی، کلاسی است که توسط اکثر درختان انتخاب شده است. شکل زیر مثالی ساده از پیاده‌سازی الگوریتم جنگل تصادفی را بر روی نمونه‌ای از میوه‌ها نشان می‌دهد.

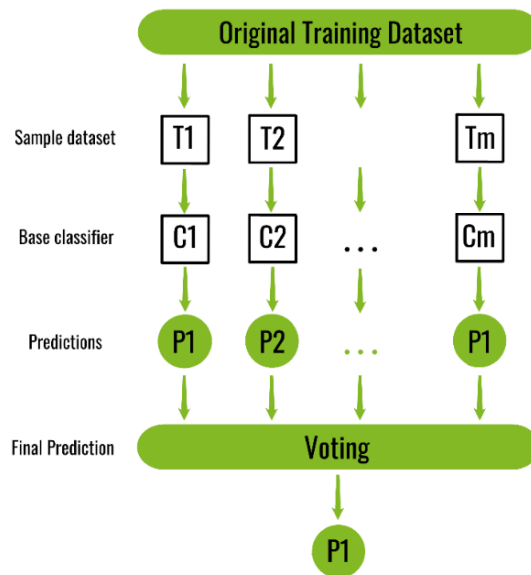


شکل (۴-۵) مثالی از پیاده‌سازی الگوریتم جنگل تصادفی برای دسته‌بندی داده‌ها

^{۶۸} Random Forest

۵-۴-۳-۵- دسته‌بندی کیسه‌ای^{۶۹}

الگوریتم دسته‌بندی کیسه‌ای یک فرا برآوردگر^{۷۰} ترکیبی^{۷۱} است که هر کدام از طبقه‌بندی‌کننده‌های پایه را بر روی زیرمجموعه‌های تصادفی مجموعه داده اصلی قرار می‌دهد و سپس پیش‌بینی‌های فردی آن‌ها (چه با رأی‌گیری^{۷۲} یا با میانگین‌گیری^{۷۳}) را جمع‌آوری می‌کند تا یک پیش‌بینی نهایی را تشکیل دهد. چنین فرا برآوردگر معمولاً می‌تواند به عنوان راهی برای کاهش واریانس تخمین‌گر جعبه سیاه^{۷۴} (به عنوان مثال، درخت تصمیم)، با ورود تصادفی به مراحل ایجاد آن و سپس ساختن مجموعه‌ای از آن استفاده شود.



شکل (۵-۵) نحوه عملکرد الگوریتم دسته‌بندی کیسه‌ای

۵-۴-۴-۵- دسته‌بندی تقویتی گرادیان^{۷۵}

طبقه‌بندی‌کننده تقویتی گرادیان یک الگوریتم یادگیری ماشین است که بسیاری از مدل‌های یادگیری ضعیف را با هم ترکیب می‌کند تا یک مدل پیش‌بینی قوی ایجاد کند. معمولاً هنگام انجام الگوریتم دسته‌بندی تقویتی گرادیان از درختان تصمیم استفاده می‌شود. مدل تقویتی گرادیان به دلیل اثربخشی در طبقه‌بندی مجموعه داده‌های پیچیده، محبوب شده‌اند.

^{۶۹} Bagging Classifier

^{۷۰} Meta-Estimator

^{۷۱} Ensemble

^{۷۲} Voting

^{۷۳} Averaging

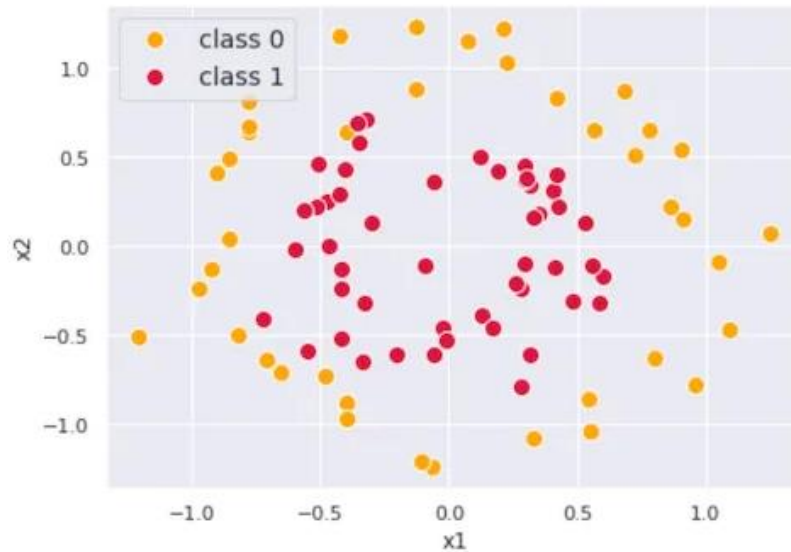
^{۷۴} Black-Box Estimator

^{۷۵} Gradient Boosting Classifier Algorithm

مدل تقویت گرادیان ترکیبی خطی از یک سری مدل‌های ضعیف است که به‌صورت تناوبی برای ایجاد یک مدل نهایی قوی ساخته شده است. این روش به خانواده الگوریتم‌های یادگیری گروهی تعلق دارد و عملکرد آن همواره از الگوریتم‌های اساسی یا ضعیف (مثلاً درخت تصمیم) یا روش‌های بر اساس ریشه‌گذاری (مانند جنگل تصادفی) بهتر است؛ اما این موضوع تا حدی از مشخصات داده‌های ورودی تأثیر می‌پذیرد. روش این الگوریتم بدین ترتیب است که تابع هزینه^{۷۶} را به کمینه‌ترین مقدار خود برساند. در علم آمار، معمولاً تابع هزینه برای اینکه مشخص شود تخمین پارامترمان تا چه حد موفق بوده، استفاده می‌شود. تابعی است که برای سنجش میزان موفقیت تخمین‌گر از تخمین پارامتر نسبت به مقادیر واقعی از آن استفاده می‌شود. در مسائل طبقه‌بندی، تابع هزینه در اصل به‌نوعی تعداد طبقه‌بندی‌های اشتباه توسط تخمین‌گر را نمایان می‌کند.

الگوریتم یادگیری تقویتی یک الگوریتم تقویتی قدرتمند است که چندین یادگیرنده ضعیف را با یادگیرندگان قوی ترکیب می‌کند که در آن هر مدل جدید برای به‌حداقل رساندن تابع هزینه مانند میانگین مربعات خطایا آنتروپی متقابل مدل قبلی، با استفاده از گرادیان نزول آموزش داده می‌شود. در هر تکرار، الگوریتم گرادیان تابع هزینه را با توجه به پیش‌بینی‌های مجموعه فعلی محاسبه می‌کند و سپس یک مدل ضعیف جدید را برای به‌حداقل رساندن این گرادیان آموزش می‌دهد. سپس پیش‌بینی‌های مدل جدید به مجموعه اضافه می‌شود و این فرایند تا زمانی که یک معیار توقف برآورده شود، تکرار می‌شود. در این الگوریتم، وزن نمونه‌های آموزشی بهینه‌سازی نشده است، در عوض، هر پیش‌بینی‌کننده با استفاده از خطاهای باقی‌مانده قبلی به‌عنوان برچسب آموزش داده می‌شود. در ادامه روند پیاده‌سازی مدل تقویتی گرادیان در قالب یک مثال آموزشی توضیح داده خواهد شد. شکل زیر نمایش داده‌های دسته‌بندی در شکل زیر نشان داده شده است.

^{۷۶} Loss Function

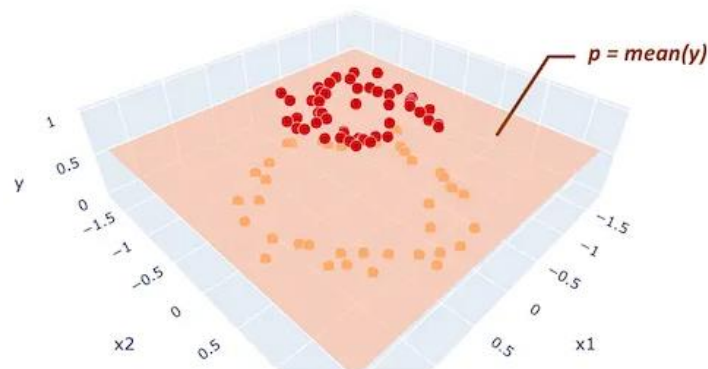


شکل (۵-۶) مثالی از مسئله دسته‌بندی دو کلاسه

هدف ساخت یک مدل تقویتی گرادیان است که داده‌ها را به دودسته دسته‌بندی کند. اولین گام، ایجاد یک پیش‌بینی یکنواخت بر روی احتمال کلاس ۱ (ما آن را p می‌نامیم) برای تمام نقاط داده است که در واقع همان میانگین کلاس می‌باشد.

$$p = P(y = 1) = \bar{y}$$

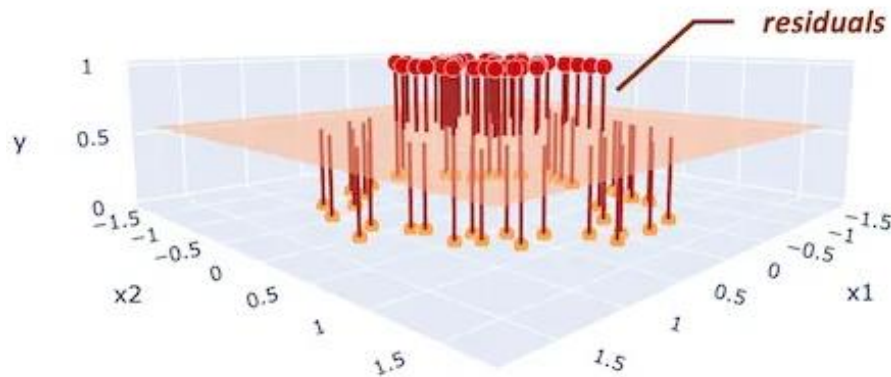
در اینجا یک نمایش سه‌بعدی از داده‌ها و پیش‌بینی اولیه آمده است. در این لحظه، پیش‌بینی فقط صفحه‌ای است که همیشه مقدار یکنواخت $p = \text{mean}(y)$ را در محور y دارد.



شکل (۵-۷) نمایش صفحه پیش‌بینی به شکل سه‌بعدی

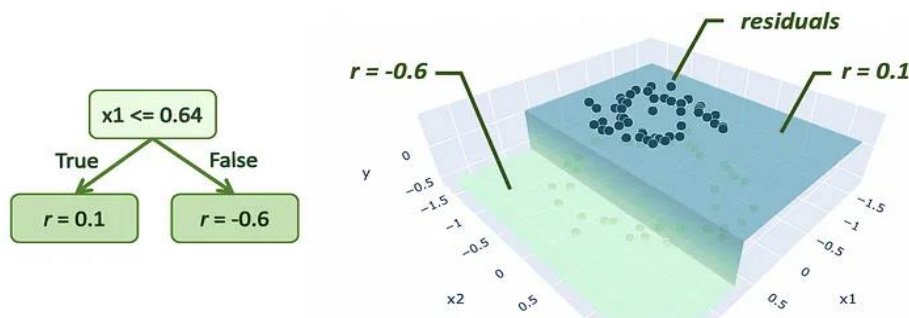
در این مثال، میانگین y ، ۰.۵۶ است. از آنجایی که بزرگ‌تر از ۰.۵ است، همه چیز با این پیش‌بینی اولیه در کلاس ۱ طبقه‌بندی می‌شود. ممکن است به نظر برسد که این پیش‌بینی ارزش یکسان، منطقی نیست، لازم به ذکر است که با اضافه کردن مدل‌های ضعیف بیشتر به آن، پیش‌بینی بهبود می‌یابد.

برای بهبود کیفیت پیش‌بینی، ممکن است روی باقی‌مانده‌ها (خطای پیش‌بینی) از پیش‌بینی اولیه تمرکز کنیم، زیرا این همان چیزی است که باید به حداقل برسد. باقی‌مانده‌ها به صورت $r_i = y_i - p$ تعریف می‌شوند (i نشان‌دهنده شاخص هر نقطه داده است). در شکل زیر باقیمانده‌ها به صورت خطوط قهوه‌ای نشان داده شده‌اند که خطوط عمود از هر نقطه داده به صفحه پیش‌بینی هستند.



شکل (۵-۸) نمایش باقی‌مانده‌ها

برای به حداقل رساندن این باقی‌مانده‌ها، یک مدل درخت رگرسیون با x_1 و x_2 به عنوان ویژگی‌های آن و باقی‌مانده r به عنوان هدف آن باید ساخته شود. اگر بتوان درختی ساخت که الگوهایی را بین x و r پیدا کند، می‌توان با استفاده از آن الگوهای یافت شده، باقی‌مانده‌های حاصل از پیش‌بینی اولیه p را کاهش داد. برای ساده کردن نمایش، درختان بسیار ساده‌ای که هر کدام فقط دارای یک تقسیم و دو گره برگ هستند، ساخته شده‌اند که به آن «استامپ»^{۷۷} می‌گویند. لازم به ذکر است که درخت‌های تقویت‌کننده گرادین معمولاً درختان کمی عمیق‌تر مانند درخت‌هایی با ۸ تا ۳۲ گره برگ دارند. در اینجا اولین درخت ایجاد شده باقی‌مانده‌ها را با دو مقدار مختلف $r = \{0.1, -0.6\}$ پیش‌بینی می‌کند.



شکل (۵-۹) درخت ساخته شده برای متغیرهای x و باقی‌مانده r

^{۷۷} Stump

در ادامه گاما طبق فرمول زیر محاسبه می‌شود. مقادیر گاما را به پیش‌بینی اولیه خود اضافه می‌کنیم تا باقی‌مانده‌ها را کاهش دهیم.

$$\gamma_j = \frac{\sum_{x_i \in R_j} (y_i - p)}{\sum_{x_i \in R_j} p(1 - p)}$$

γ is computed for each terminal node j

Aggregating for all the data points x_i that belongs to terminal node j

مقادیر گاما ۱ و ۲ بدین ترتیب محاسبه می‌شوند.

$$\gamma_1 = \frac{\sum_{x_i \in R_1} (y_i - 0.56)}{\sum_{x_i \in R_1} 0.56 \cdot (1 - 0.56)} = 0.3$$

$$\gamma_2 = \frac{\sum_{x_i \in R_2} (y_i - 0.56)}{\sum_{x_i \in R_2} 0.56 \cdot (1 - 0.56)} = -2.2$$

برای اینکه گاما را به مقدار p اضافه شود، به شکل زیر عمل می‌کنیم. ابتدا مقدار $\log(\text{odds})$ را از p به دست می‌آوریم (به آن $F(x)$ گفته می‌شود). سپس گاما را به آن اضافه می‌کنیم.

$$\log(\text{odds}) = \log\left(\frac{p}{1-p}\right)$$

برای اینکه مدل بیش از حد از آموزش^{۷۸} نبیند و خطای آن کاهش یابد می‌توان مقدار گاما را با یک‌وزنی (بین ۰ تا ۱) که به آن نرخ یادگیری^{۷۹} v گفته می‌شود، ضرب کرد و سپس به مقدار $\log(\text{odds})$ یا همان $F(x)$ اضافه نمود تا پیش‌بینی بروز شود.

$$F_1(x) = F_0(x) + v \cdot \gamma$$

Updated prediction

Initial prediction

Learning rate

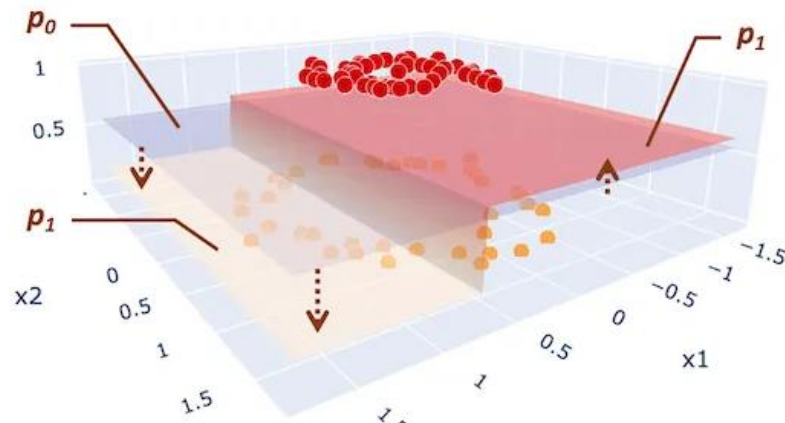
در این مثال، ما از نرخ یادگیری نسبتاً بزرگ $v = 0.9$ استفاده می‌کنیم تا فرایند بهینه‌سازی را آسان‌تر درک کنیم، اما معمولاً قرار است مقادیر بسیار کوچک‌تری مانند ۰.۱ در نظر گرفته شود. با جایگزینی مقادیر واقعی برای متغیرهای سمت راست معادله بالا، پیش‌بینی به‌روز $F_1(x)$ را به دست می‌آید.

^{۷۸} Overfit

^{۷۹} Learning Rate

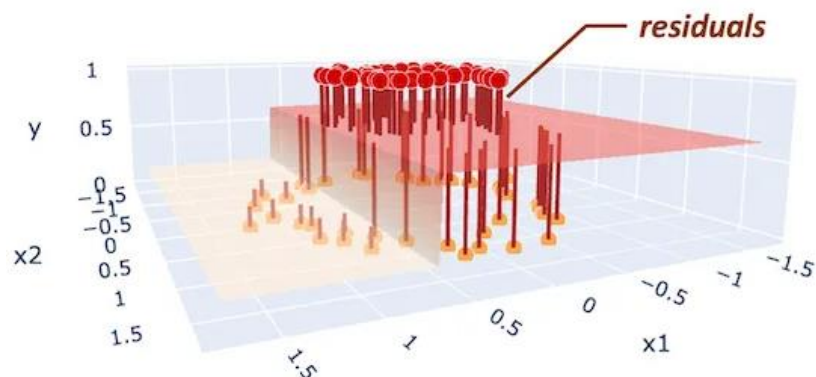
$$F_1(x) = \begin{cases} \log\left(\frac{0.56}{1-0.56}\right) + 0.9 \cdot 0.3 = 0.5 & \text{if } x_1 \leq 0.64 \\ \log\left(\frac{0.56}{1-0.56}\right) - 0.9 \cdot 2.2 = -1.7 & \text{otherwise} \end{cases}$$

اگر $\log(\text{odds})$ را دوباره به p تبدیل کنیم. شکی پله‌مانند از داده‌ها مانند شکل زیر به دست می‌آید.



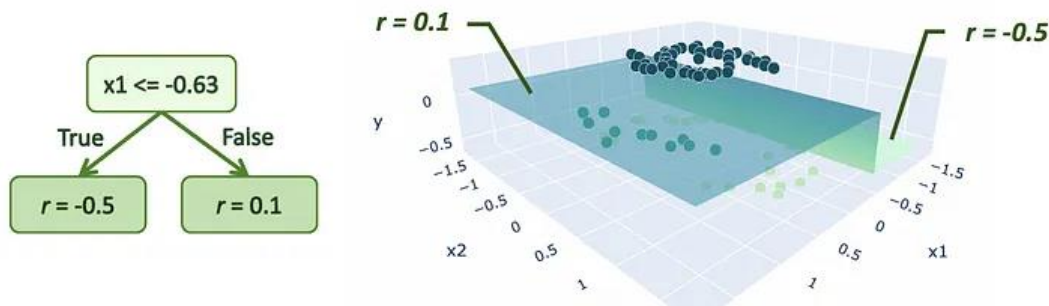
شکل (۵-۱۰) نمایش پیش‌بینی بروز شده

اکنون، باقی‌مانده‌های بروز شده r به شکل زیر است.



شکل (۵-۱۱) نمایش باقی‌مانده‌های بروز شده

مجدداً یک درخت رگرسیون با استفاده از همان x_1 و x_2 به‌عنوان ویژگی‌های ورودی برای باقی‌مانده‌های بروز شده ایجاد می‌کنیم.



شکل (۵-۱۲) درخت ساخته شده برای متغیرهای x و باقی مانده r بروز شده

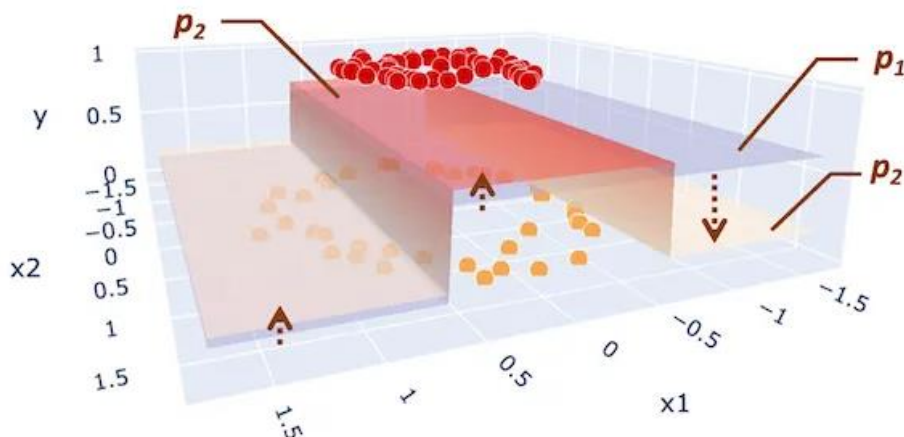
حال مجدداً گاما به همان روش قبل محاسبه کرده و $F_2(x)$ را به دست می آوریم.

$$F_2(x) = \begin{cases} F_1(x) - v \cdot 2.3 = 0.5 - 0.9 \cdot 2.3 = -1.6 & \text{if } x_1 \leq -0.63 \\ F_1(x) + v \cdot 0.4 = 0.5 + 0.9 \cdot 0.4 = 0.9 & \text{else if } -0.63 < x_1 \leq 0.64 \\ F_1(x) + v \cdot 0.4 = -1.7 + 0.9 \cdot 0.4 = -1.3 & \text{otherwise} \end{cases}$$

These are γ computed with this formula:

$$\gamma_j = \frac{\sum_{x_i \in R_j} (y_i - p)}{\sum_{x_i \in R_j} p(1 - p)}$$

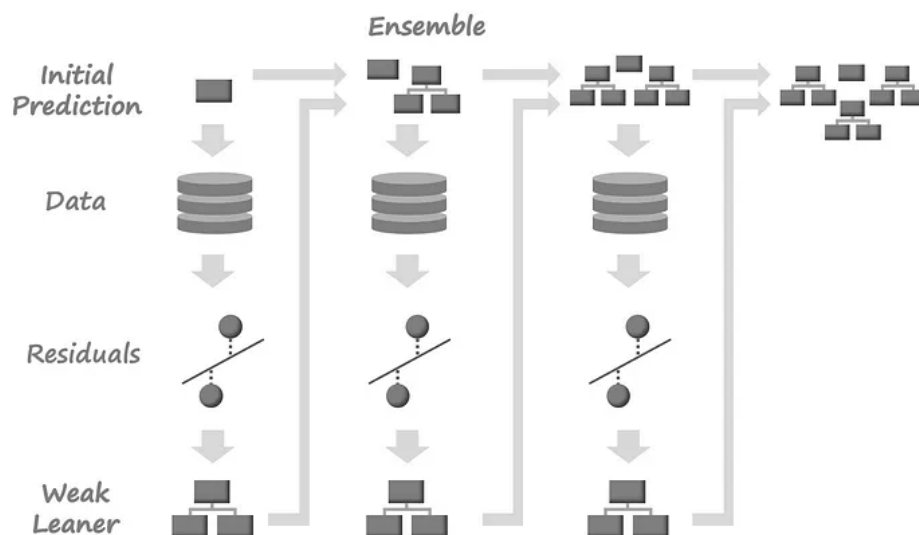
سپس $F_2(x)$ را به $p_2(x)$ تبدیل می کنیم و به شکل زیر می رسمیم.



شکل (۵-۱۳) نمایش پیش بینی بروز شده

سپس، این مراحل را تکرار می کنیم تا زمانی که پیش بینی مدل متوقف شود. می توان دید که پیش بینی ترکیبی $p(x)$ به هدف ما نزدیک تر می شود؛ زیرا درخت های بیشتری را به مدل ترکیبی اضافه می کنیم. این روشی است که الگوریتم تقویتی گرادیان برای پیش بینی اهداف پیچیده با ترکیب چندین مدل ضعیف انجام می دهد.

تصویر زیر به طور خلاصه کل فرایند این الگوریتم را نشان می دهد.



شکل (۵-۱۴) فرایند الگوریتم تقویتی گرادیان

۵-۴-۴-۵- دسته‌بندی XGboost^{۸۰}

XGBoost یک الگوریتم یادگیری ماشین مبتنی بر درخت تصمیم است که از یک چارچوب تقویت گرادیان استفاده می‌کند. شکل زیر فرایند بهینه‌سازی الگوریتم ماشین تقویت گرادیان^{۸۱} را توسط XGboost نشان می‌دهد.



شکل (۵-۱۵) نحوه بهینه‌سازی در الگوریتم XGboost

^{۸۰} XGboost Classifier

^{۸۱} Gradient Boosting Machines (GBMs)

XGBoost مخفف واژه تقویت گرادیان شدید^{۸۲} است و به دلیل توانایی آن در مدیریت مجموعه داده‌های بزرگ و توانایی آن برای دستیابی به عملکرد پیشرفته در بسیاری از وظایف یادگیری ماشین، به یکی از محبوب‌ترین و پرکاربردترین الگوریتم‌های یادگیری ماشین تبدیل شده است.

در این الگوریتم درخت‌های تصمیم به صورت متوالی ایجاد می‌شوند. وزن‌ها نقش مهمی در XGBoost دارند. وزن‌ها به همه متغیرهای مستقل اختصاص داده می‌شوند که سپس به درخت تصمیم که نتایج را پیش‌بینی می‌کند، وارد می‌شوند. وزن متغیرهای پیش‌بینی شده اشتباه توسط درخت، افزایش می‌یابد و این متغیرها سپس به درخت تصمیم دوم تغذیه می‌شوند. سپس این طبقه‌بندی‌کننده‌ها یا پیش‌بینی‌کننده‌های منفرد برای ارائه یک مدل قوی و دقیق‌تر جمع می‌شوند.

۵-۵- اعتبارسنجی و بررسی صحت الگوریتم‌های رگرسیون و دسته‌بندی

پس از پیاده‌سازی و آموزش مدل‌ها روی داده‌های آموزشی، اعتبارسنجی الگوریتم‌ها با داده‌های اعتبارسنجی انجام شد. در ادامه به معیارهای اعتبارسنجی الگوریتم‌ها می‌پردازیم.

۵-۵-۱- میانگین مربعات خطا^{۸۳}

یکی از معیارهای ارزیابی الگوریتم‌های رگرسیون می‌باشد. فرمول آن به صورت زیر است.

$$MSE(Y, \hat{Y}) = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \frac{1}{n} \sum_{i=1}^n e_i^2$$

در واقع در این معیار ارزیابی، خطاها به توان دو رسیده، با هم جمع می‌شوند و بر تعداد نقاط تقسیم می‌شوند. نکته‌ای که در مورد این معیار وجود دارد، وزن‌دهی بیشتر به خطاهای بزرگ‌تر است، به طوری که اگر دو داده با خطاهای ۱ و ۳ وجود داشته باشد، اثرگذاری داده دوم بیشتر خواهد بود.

۵-۵-۲- جذر میانگین مربعات خطا^{۸۴}

معیار جذر میانگین مربعات خطا، یک معیار ارزیابی رگرسیون است. با استفاده از رابطه زیر از میانگین مربعات خطا محاسبه می‌شود.

$$RMSE(Y, \hat{Y}) = \sqrt{MSE(Y, \hat{Y})}$$

دلیل استفاده از جذر میانگین مربعات خطا، یکسان بودن بُعد و مقیاس آن با ویژگی هدف است. برای مثال، اگر یک مدل برای پیش‌بینی وزن افراد برحسب Kg ایجاد کرده باشیم، واحد میانگین مربعات خطا

^{۸۲} Extreme Gradient Boosting

^{۸۳} MSE: Mean Squared Error

^{۸۴} RMSE: Root Mean Squared Error

R^2 Kg خواهد بود، درحالی که که واحد جذر مربعات خطا Kg است. این معیار اغلب در گزارش نتایج استفاده می شود و به عنوان تابع هزینه استفاده نمی شود.

۵-۳-۵- میانگین قدرمطلق خطا^{۸۵}

در این معیار ارزیابی رگرسیون، به جای به توان ۲ رساندن خطاها، از تابع قدرمطلق استفاده می شود و به صورت زیر محاسبه می شود.

$$MAE(Y, \hat{Y}) = \frac{1}{n} \sum_{i=1}^n |Y_i - \hat{Y}_i| = \frac{1}{n} \sum_{i=1}^n |e_i|$$

۵-۴-۵- ضریب امتیاز R^2

این معیار برخلاف ۳ معیار قبلی، با افزایش، دقت بالای مدل را نشان می دهد. برای محاسبه آن به شکل زیر عمل می کنیم.

$$R^2(Y, \hat{Y}) = 1 - \frac{\sum_{i=1}^n (Y_i - \hat{Y}_i)^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2} = 1 - \frac{MSE(Y, \hat{Y})}{\sigma^2(Y)}$$

این معیار همواره عددی کوچک تر از ۱ است. اگر مدلی همواره میانگین ویژگی هدف را در خروجی تولید کند، مقدار این معیار برابر ۰ خواهد بود. ضریب تعیین معمولاً به صورت درصد بیان می شود. از ضریب تعیین تنها برای مقایسه مدل ها و گزارش نتایج استفاده می شود.

۵-۵-۵- مربع R تنظیم شده^{۸۶}

این معیار ارزیابی رگرسیون، نوع بهبودیافته معیار ارزیابی مربع R به حساب می آید. مشکلی که روش قبل دارد این است که با افزایش ویژگی های مدل، مقدار R^2 نیز افزایش می یابد که البته باعث به وجود آمدن مدل خوبی نیز می شود. اما مربع R تنظیم شده برای حل این مشکل ارائه شده است. این روش فقط ویژگی هایی را در نظر می گیرد که برای مدل مهم هستند و براین اساس می تواند بهبود واقعی مدل را نشان دهد. همچنین، معیار R تنظیم شده همیشه از معیار R^2 کمتر است. معادله این معیار به شکل زیر است.

$$R^2_{adjusted} = 1 - \frac{(1 - R^2)(N - 1)}{N - P - 1}$$

در رابطه فوق، R^2 نشان دهنده تعداد نمونه های مربع R، آرگومان p نشان دهنده تعداد پیش بینی ها و N اندازه کل نمونه ها هستند.

^{۸۵} MAE: Mean Absolute Error

^{۸۶} Adjusted R squared

۵-۵-۶- ماتریس اغتشاش^{۸۷}

برای اینکه بتوانیم نتایج دسته‌بندی الگوریتم را با داده‌های واقعی مقایسه کنیم، از ماتریس اغتشاش استفاده می‌کنیم. جدول زیر، ماتریس اغتشاش را نشان می‌دهد.

جدول (۵-۲) ماتریس اغتشاش

	Predicted Positive	Predicted Negative
Actual Positive	True Positive (TP)	False Negative (FN)
Actual Negative	False Positive (FP)	True Negative (TN)

سطرها مقادیر واقعی و ستون‌ها مقادیر پیش‌بینی شده را نشان می‌دهند. سلول‌های این ماتریس مفاهیم زیر را ارائه می‌دهند.

- **مثبت - صحیح^{۸۸}**: نشان می‌دهد که مدل یک نتیجه مثبت را پیش‌بینی کرده است و مشاهده واقعی درست بوده است.
- **مثبت - کاذب^{۸۹}**: نشان می‌دهد که مدل یک نتیجه مثبت را پیش‌بینی کرده است، اما مشاهده واقعی نادرست بوده است.
- **منفی - کاذب^{۹۰}**: نشان می‌دهد که مدل یک نتیجه منفی را پیش‌بینی کرده است، درحالی‌که مشاهده واقعی درست بوده است.
- **منفی - صحیح^{۹۱}**: نشان می‌دهد که مدل یک نتیجه منفی را پیش‌بینی کرده است، و نتیجه واقعی نیز نادرست بوده است.

۵-۵-۷- دقت^{۹۲}

دقت معمولاً برای قضاوت در مورد عملکرد مدل استفاده می‌شود، فرمول محاسبه دقت به شکل زیر است.

^{۸۷} Confusion Matrices

^{۸۸} True Positive

^{۸۹} False Positive

^{۹۰} False Negative

^{۹۱} True Negative

^{۹۲} Accuracy

$$Accuracy = \frac{TP + TN}{TP + FP + FN + TN}$$

در واقع دقت، میزان پیش‌بینی درست مدل را بر کل محاسبه می‌کند.

۵-۸-۵- صحت^{۹۳}

این معیار، اندازه‌گیری مثبت‌های واقعی نسبت به تعداد کل مثبت‌های پیش‌بینی شده توسط مدل را محاسبه می‌کند. در واقع این معیار، میزان مثبت‌بودن پیش‌بینی‌های مثبت مدل را اندازه‌گیری می‌کند.

$$Precision = \frac{TP}{TP + FP}$$

۵-۹-۵- پوشش^{۹۴}

معیار پوشش قادر به سنجش مثبت پیش‌بینی شده مدل نسبت به تعداد پیامدهای مثبت واقعی است. با استفاده از این معیار، می‌توان ارزیابی کرد که مدل چقدر قادر به شناسایی نتایج واقعی است.

$$Recall = \frac{TP}{TP + FN}$$

۵-۱۰-۵- امتیاز F1^{۹۵}

این معیار، میانگین هارمونیک بین دقت و پوشش است. فرمول آن به شکل زیر است.

$$F1\ Score = \frac{2 * (Precision * Recall)}{Precision + Recall}$$

این معیار هنگامی که داده‌ها به صورت نامتوازن پخش شده‌اند، دید بهتری از عملکرد مدل ارائه می‌دهد.

۵-۱۱-۵- اعتبارسنجی متقابل^{۹۶}

اعتبارسنجی متقابل یک روش آماری است که برای تخمین عملکرد مدل‌های یادگیری ماشین استفاده می‌شود. این روش برای ارزیابی چگونگی تعمیم نتایج یک تحلیل آماری به یک مجموعه داده دیده نشده است. این روش، آموزش بیش از حد مدل را شناسایی می‌کند و با بررسی داده‌های دیده نشده نتیجه دقیق‌تری راجع به عملکرد مدل، ارائه می‌دهد.

در این پژوهش ما از روش اعتبارسنجی متقابل مونت کارلو^{۹۷} برای ارزیابی نهایی عملکرد الگوریتم‌ها روی داده‌های دیده نشده، استفاده کرده‌ایم. در ادامه به نحوه عملکرد این معیار اعتبارسنجی می‌پردازیم.

این معیار، یک استراتژی بسیار انعطاف‌پذیر برای اعتبارسنجی متقابل است. در این تکنیک، مجموعه داده‌ها به طور تصادفی به مجموعه‌های آموزشی و اعتبارسنجی تقسیم می‌شوند. درصدی از مجموعه داده‌ای را که قرار است به عنوان مجموعه آموزشی استفاده شود و درصدی که به عنوان مجموعه اعتبارسنجی

^{۹۳} Precision

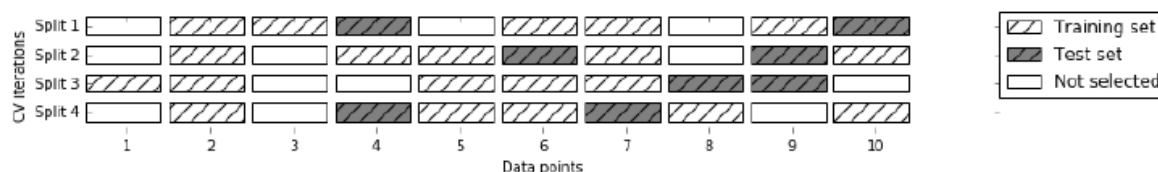
^{۹۴} Recall

^{۹۵} F1 Score

^{۹۶} Cross Validation

^{۹۷} Monte Carlo Cross-Validation (Shuffle Split)

استفاده می‌شود، را مشخص می‌کنیم. اگر مجموع درصدها به ۱۰۰ نرسد، از مجموعه داده باقی مانده استفاده نمی‌شود. سپس این تقسیم‌بندی به تعداد دفعاتی که مشخص می‌کنیم، تکرار می‌شود و دقت هر تکرار محاسبه می‌گردد. می‌توان میانگین نهایی دقت دفعات تکرار را به عنوان میزان دقت نهایی مدل روی داده‌های دیده نشده در نظر گرفت.



شکل (۵-۱۶) نحوه عملکرد اعتبارسنجی متقابل مونت کارلو

۵-۵-۱۲- ROC - AUC^{۹۸}

این معیار ارزیابی روی داده‌های دیده نشده برای مقایسه نهایی مدل‌ها، پیاده‌سازی شده است. برای درک بهتر این معیار، مفاهیم زیر مطرح می‌گردد.

- **نرخ مثبت صحیح^{۹۹}:** یک معیار ارزیابی عملکرد می‌باشد. همان مفهوم معیار پوشش می‌باشد و مشخص می‌کند که به چه نسبتی پیش‌بینی صحیح صورت گرفته است. فرمول آن به شکل زیر است.

$$TPR = \frac{TP}{P} = \frac{TP}{TP + FN}$$

- **نرخ مثبت کاذب^{۱۰۰}:** یک معیار ارزیابی عملکرد می‌باشد و نشانگر تعداد شناسایی‌های مثبت از میان مشاهدات منفی است. فرمول آن به شکل زیر است.

$$FPR = \frac{FP}{N} = \frac{FP}{FP + TN}$$

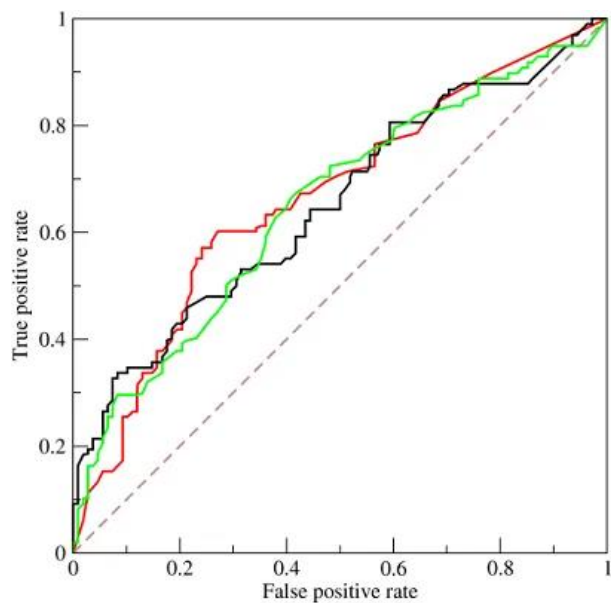
- **منحنی مشخصه عملکرد^{۱۰۱}:** یک منحنی مشخصه عملکرد، یک نمودار برای نمایش توانایی ارزیابی یک سیستم دسته‌بندی دودویی محسوب می‌شود که آستانه تشخیص آن نیز متغیر است. این منحنی توسط ترسیم نرخ مثبت صحیح بر حسب نرخ مثبت کاذب، ایجاد می‌شود. نمودار زیر منحنی مشخصه عملکرد را برای سه مدل دسته‌بندی مختلف نشان می‌دهد.

^{۹۸} AUC (Area Under the Curve) - ROC (Receiver Operating Characteristics) curve

^{۹۹} True Positive Rate (TPR)

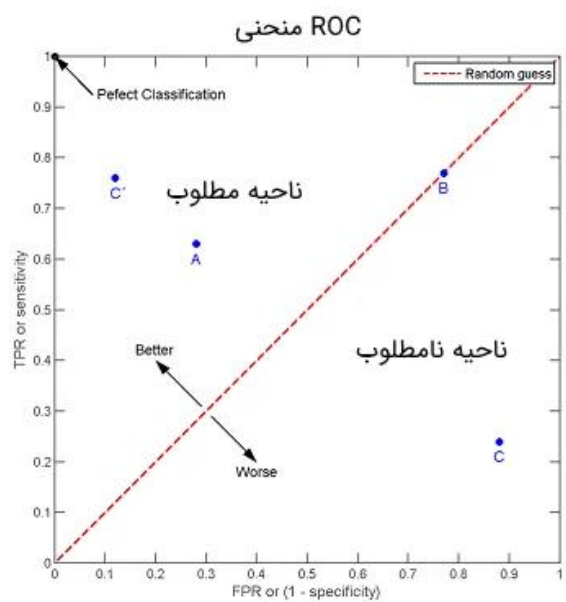
^{۱۰۰} False Positive Rate (FPR)

^{۱۰۱} ROC: Receiver Operating Characteristics curve



نمودار (۵-۱) منحنی مشخصه عملکرد برای سه روش مختلف دسته‌بندی

باتوجه به نمودار زیر، بهترین عملکرد دسته‌بندی در این نمودار در نقطه‌ای با مختصات $(0, 1)$ رخ خواهد داد که در آن کمترین نرخ اشتباه و بیشترین نرخ بازیابی یا حساسیت را داریم. این نقطه بیانگر «بهترین دسته‌بندی»^{۱۰۲} است.



نمودار (۵-۲): نواحی مطلوب و نامطلوب در منحنی ROC

^{۱۰۲} Perfect Classification

همچنین در نمودار فوق، خط منقطعی که از میان نمودار عبور کرده و نقطه $(0,0)$ را به $(1,1)$ پیوند می‌دهند، حدس تصادفی است که به صورت ناحیه $50\%-50\%$ نیز شناخته می‌شود. اگر نقطه‌ای روی این خط منقطع قرار گرفته باشد، تشخیص درستی نسبت به قرارگیری در هر گروه، برایش وجود ندارد. در حقیقت در نیمی از موارد می‌تواند در یک دسته و در نیمی از موارد نیز در دسته دیگری، طبقه‌بندی شود و نقشی در تعیین خطا نخواهد داشت. یکی از نمونه‌های معروف برای دسته‌بندی به صورت تصادفی، تصمیم تعلق نقطه به هر یک از دو گروه به وسیله پرتاب سکه است. هر چه تعداد نمونه‌ها در دسته‌بندی تصادفی بیشتر شود، این خط به قطر نواحی ROC نزدیک‌تر خواهد شد.

AUC: مساحت زیر منحنی مشخصه عملکرد را AUC می‌گویند که نشان‌دهنده درجه یا معیار تفکیک‌پذیری است. این معیار نشان می‌دهد که مدل چقدر می‌تواند بین کلاس‌ها تمایز قائل شود. یک مدل عالی دارای $AUC = 1$ است

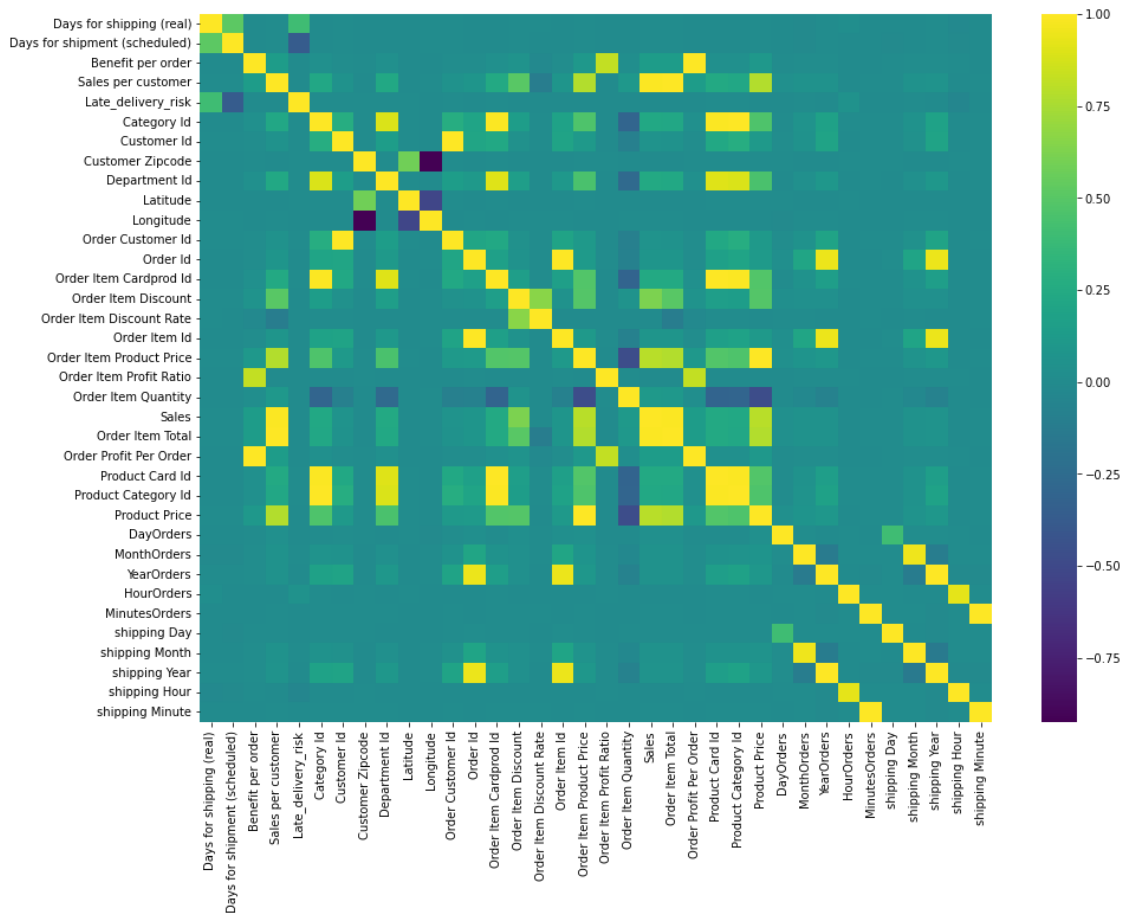
۵-۶- بررسی نتایج و مقایسه الگوریتم‌ها

در این بخش، مقایسه الگوریتم‌ها بر اساس معیارهای اعتبارسنجی، صورت گرفت و نتایج در قالب جدول ارائه گردیده است.

۶- یافته‌های تحقیق

۶-۱- تجزیه و تحلیل اکتشافی

ابتدا همبستگی میان ویژگی‌ها بررسی شده است. شکل زیر، همبستگی بین ویژگی‌ها را نشان می‌دهد. همبستگی یک معیار آماری است که میزان ارتباط خطی دو متغیر را بیان می‌کند (به این معنی که آن‌ها با هم با یک نرخ ثابت تغییر می‌کنند). این یک ابزار رایج برای توصیف روابط ساده بدون اظهارنظر در مورد علت و معلول است. همبستگی، عددی بین ۱- و ۱+ است. همبستگی مثبت، نشان‌دهنده میزان افزایش یا کاهش آن متغیرها به صورت موازی است و همبستگی منفی نشان‌دهنده میزان افزایش یک متغیر با کاهش متغیر دیگر است.



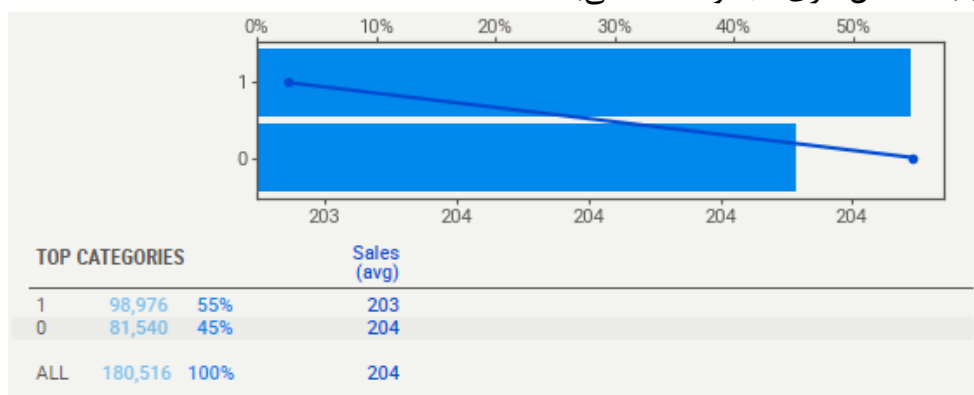
شکل (۶-۱) همبستگی ویژگی‌ها در مجموعه داده

همان‌طور که از شکل فوق مشاهده می‌گردد، ویژگی‌های مجموعه فروش انجام شده به‌ازای هر مشتری، مقدار تخفیف کالای سفارشی، قیمت محصولات بدون تخفیف، مبلغ کل به‌ازای هر سفارش و قیمت محصول

با میزان فروش همبستگی بالایی (بین ۰.۷۵ تا ۱) دارند. همچنین ویژگی روزهای ارسال واقعی محصول خریداری شده با ویژگی‌های روزهای تحویل برنامه‌ریزی شده محصول خریداری شده و وجود ریسک تأخیر در ارسال همبستگی نسبتاً بالایی دارد (بین ۰.۵ تا ۰.۷۵).

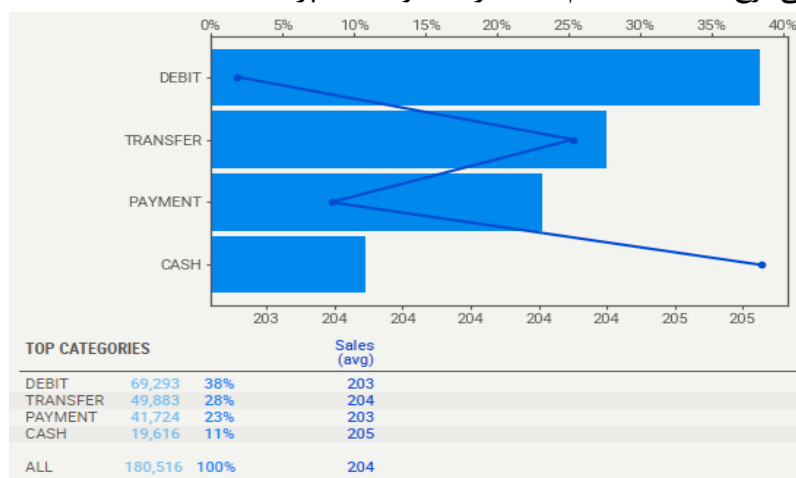
ویژگی درآمد به‌ازای هر سفارش ثبت شده، با ویژگی‌های نرخ سود کالای سفارشی و سود سفارش در هر سفارش همبستگی بالایی (بین ۰.۷۵ تا ۱) دارد. ویژگی مجموع فروش انجام شده به‌ازای هر مشتری، با ویژگی‌های قیمت محصولات بدون تخفیف، میزان فروش و مبلغ کل به‌ازای هر سفارش همبستگی بالایی (بین ۰.۷۵ تا ۱) دارد. همبستگی سایر ویژگی‌ها نیز در شکل فوق نشان داده شده است.

متغیر هدف دوم در این مطالعه، وجود ریسک ارسال با تأخیر می‌باشد. همان‌طور که در شکل زیر مشاهده می‌شود، در مجموعه داده مورد بررسی، ۵۵ درصد سفارش‌ها ریسک ارسال با تأخیر را دارند و ۴۵ درصد سفارش‌ها ریسک ارسال با تأخیر ندارند. این موضوع نشان می‌دهد که مجموعه داده مورد بررسی متعادل بوده و نیازی به متعادل‌سازی مجموعه داده نمی‌باشد.



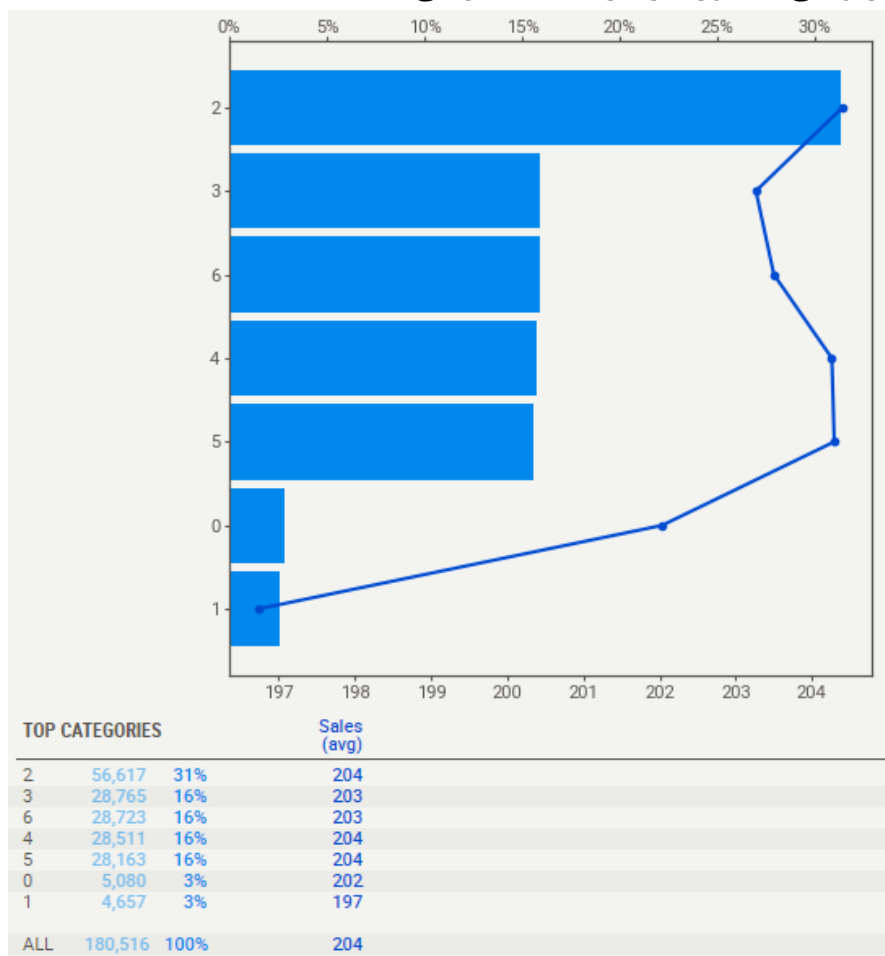
نمودار (۱-۶) متغیر هدف وجود ریسک ارسال با تأخیر

در ادامه به بررسی نوع معاملات انجام شده در مجموعه داده پرداخته شده است.



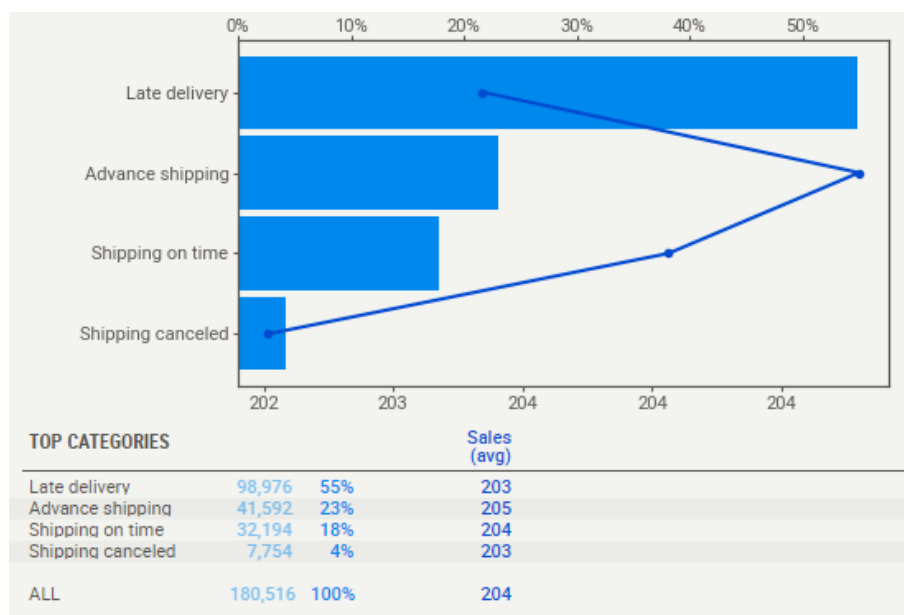
نمودار (۲-۶) نمودار انواع معاملات انجام شده

همان‌طور که مشاهده می‌شود، بیشترین میزان معاملات در این مجموعه داده، مربوط پرداخت وکالتی (۶۹۲۹۳ سفارش معادل ۳۸ درصد سفارش‌ها) می‌باشد. سفارش‌ها در حال پرداخت (۲۸ درصد معادل ۴۹۸۸۳ سفارش)، سفارش‌ها پرداخت نشده (۲۳ درصد معادل ۴۱۷۲۴ سفارش) و پرداخت آنی (۱۱ درصد معادل ۱۹۶۱۶ سفارش) رتبه‌های دوم تا چهارم را در میان انواع معاملات انجام شده در اختیار دارند. نمودار زیر روزهای ارسال واقعی محصول خریداری شده نشان می‌دهد.



نمودار (۳-۶) نمودار روزهای ارسال واقعی محصول خریداری شده

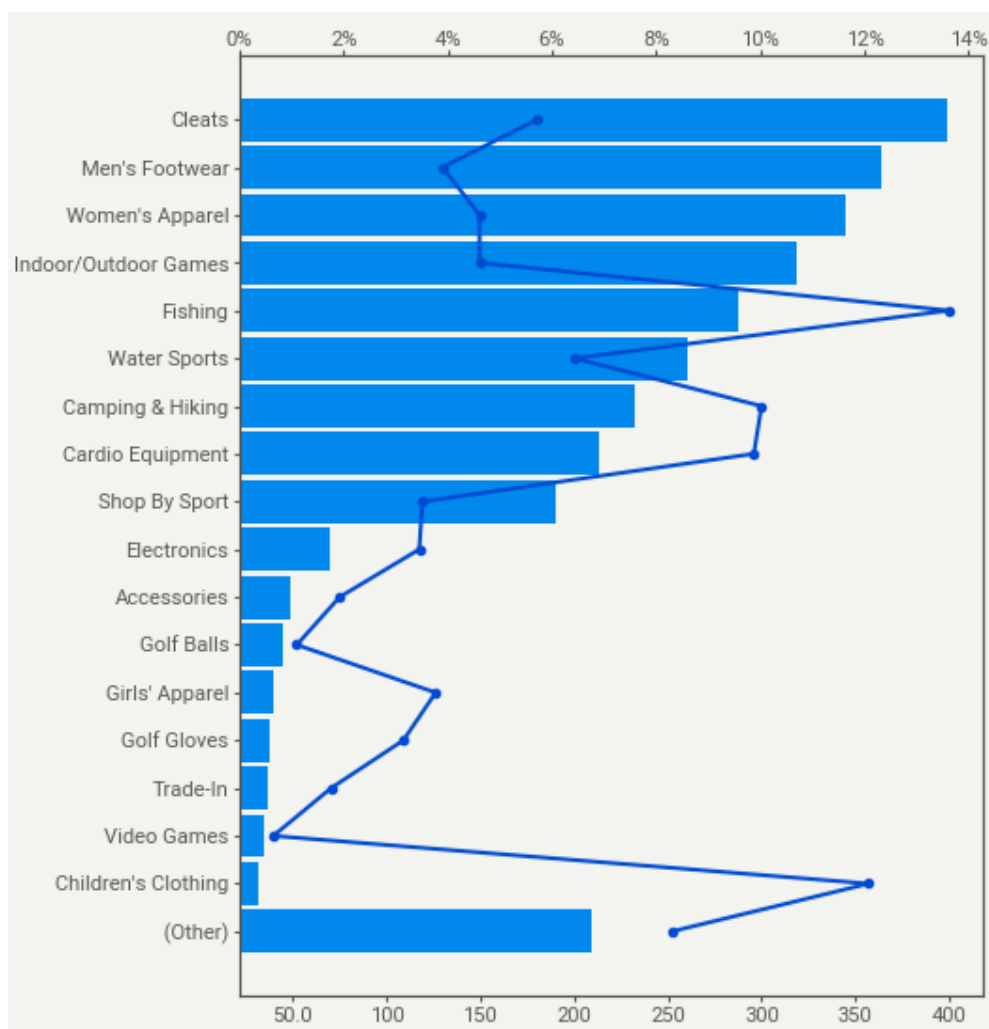
همان‌طور که از نمودار فوق قابل مشاهده است، ۳۱ درصد سفارش‌ها در ۲ روز به مشتری تحویل داده شده‌اند. سایر روزهای ارسال واقعی محصول خریداری شده در نمودار فوق نشان داده شده است. در ادامه به توزیع ویژگی وضعیت تحویل سفارش‌ها در مجموعه داده مورد بررسی می‌پردازیم.



نمودار (۴-۶) نمودار وضعیت تحویل سفارش‌ها در مجموعه داده

همان‌طور که در نمودار فوق مشاهده می‌گردد، ۵۵ درصد سفارش‌ها تأخیر در زمان تحویل دارند و در زمان برنامه‌ریزی شده به مشتری تحویل داده نشده‌اند. ۲۳ درصد سفارش‌ها پیش از زمان مقرر و برنامه‌ریزی شده به مشتری تحویل داده شده‌اند. همچنین ۱۸ درصد سفارش‌ها در زمان مقرر به مشتری تحویل داده شده است و تنها ۴ درصد سفارش‌ها توسط مشتری لغو شده‌اند که باید علل ریشه‌ای لغو این سفارش‌ها مورد بررسی قرار گیرد تا برای اجرای پروژه‌های بهبود، اصلاحات لازم صورت بگیرد. در اینجا تمام سفارشات لغو شده‌اند، سفارش‌ها در حال پرداخت بوده‌اند. ممکن است یکی از علل لغو سفارش، اختلال در درگاه پرداخت فروشگاه، بوده باشد.

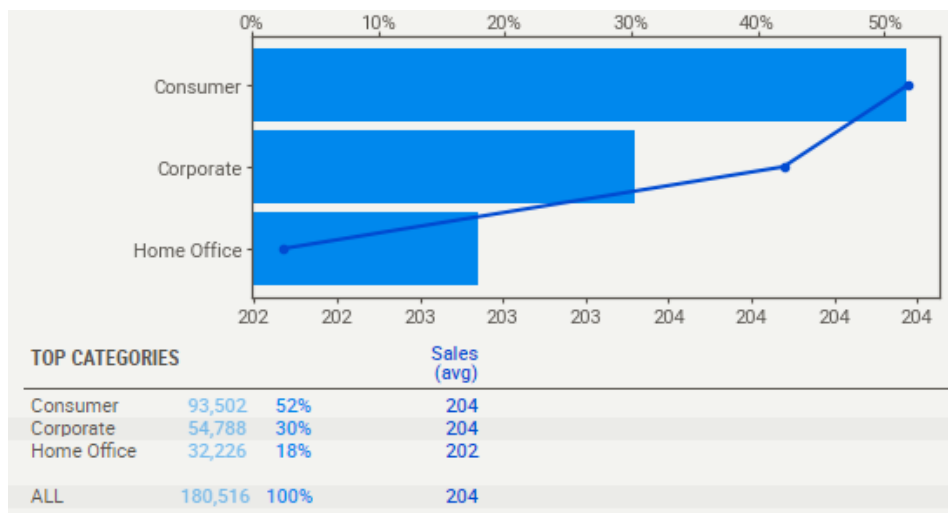
در ادامه به بررسی دسته‌بندی محصولات می‌پردازیم. نمودار فراوانی دسته‌بندی محصولات در ادامه آمده است.



نمودار (۵-۶) نمودار فراوانی دسته‌بندی محصولات سفارش داده شده

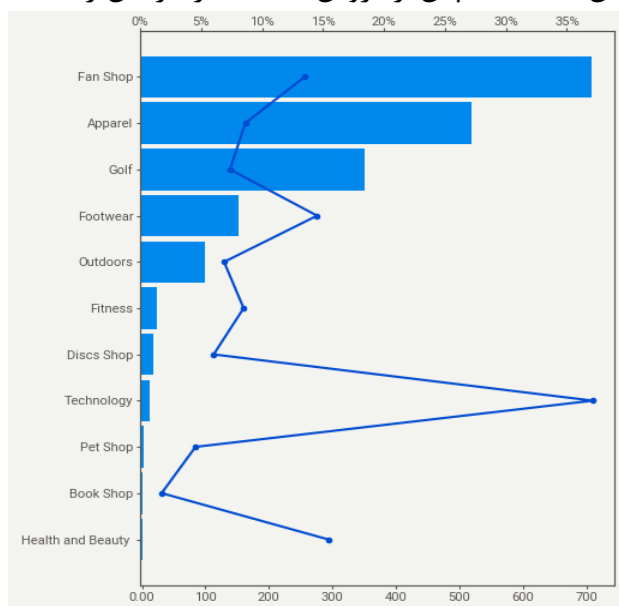
باتوجه به نمودار فوق، می‌توان متوجه شد که ۱۴ درصد محصولات سفارش داده شده کفش ورزشی می‌باشد. ۱۲ درصد محصولات سفارش داده شده کفش مردانه و ۱۲ درصد محصولات نیز پوشاک زنانه هستند. همچنین ۱۱ درصد محصولات سفارش داده شده مربوط به بازی‌های داخل و خارج از خانه هستند. ماهیگیری در رتبه بعدی دسته‌بندی محصولات سفارش داده است. سایر محصولات به ترتیب مربوط به ورزش‌های آبی، تجهیزات کمپینگ و پیاده‌روی، تجهیزات ورزشی و غیره هستند. باتوجه به تعداد سفارش‌ها انجام شده می‌توان متوجه شده که محصولات کفش ورزشی، کفش مردانه و پوشاک زنانه نزدیک به ۴۰ درصد سفارش‌ها فروشگاه را شامل می‌شوند که فروشگاه باتوجه به این موضوع باید برای فراهم کردن این محصولات برنامه‌ریزی دقیق‌تری بر اساس سلیقه مصرف‌کننده داشته باشد.

نمودار بخش‌بندی مشتریان فروشگاه در زیر آمده است.



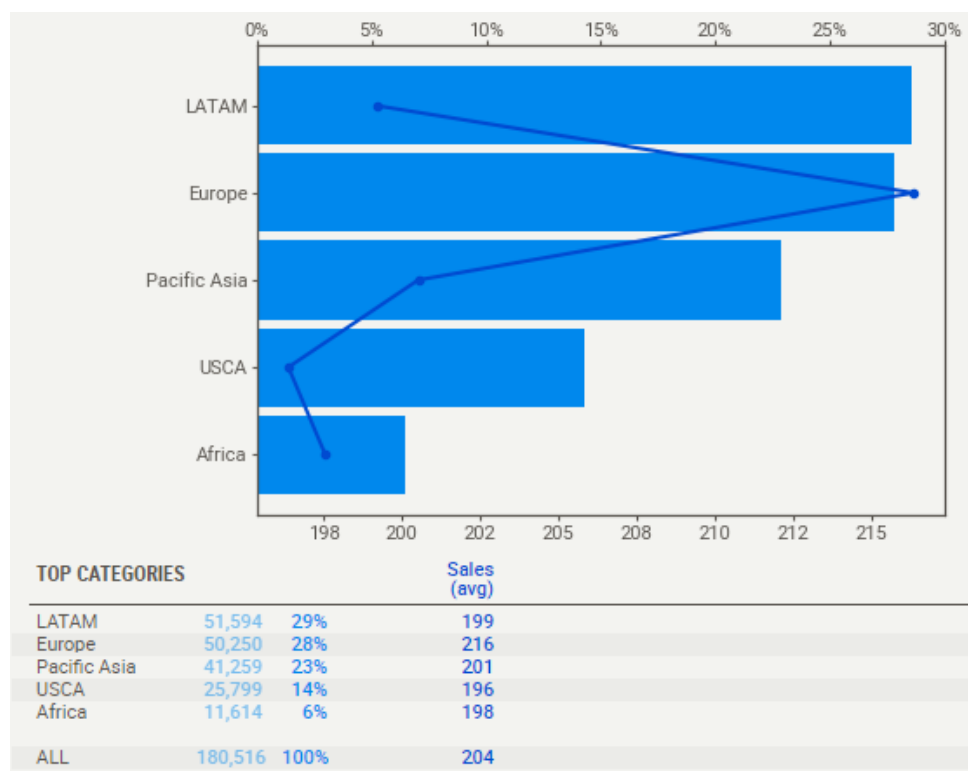
نمودار (۶-۶) نمودار بخش‌بندی مشتریان فروشگاه

باتوجه به نمودار فوق، ۵۲ درصد مشتریان فروشگاه مصرف‌کنندگان هستند. ۳۰ درصد مشتریان، شرکت‌های بزرگ و ۱۸ درصد مشتریان فروشگاه، شرکت‌های کوچک هستند. توجه به هر کدام از این مشتریان از اهمیت خاص خود برخوردار است. از آنجایی که میانگین فروش این مشتریان تقریباً با یکدیگر برابر است، بنابراین فروشگاه باید برای جلب نظر تمام بخش‌های مشتریان تلاش کند و محصولات درخواستی آن‌ها را فراهم آورد. همچنین با خدمات پس از فروش مناسب در افزایش رضایت مشتریان تلاش نماید.



نمودار (۶-۷) نمودار فراوانی دپارتمان‌های فروشگاه

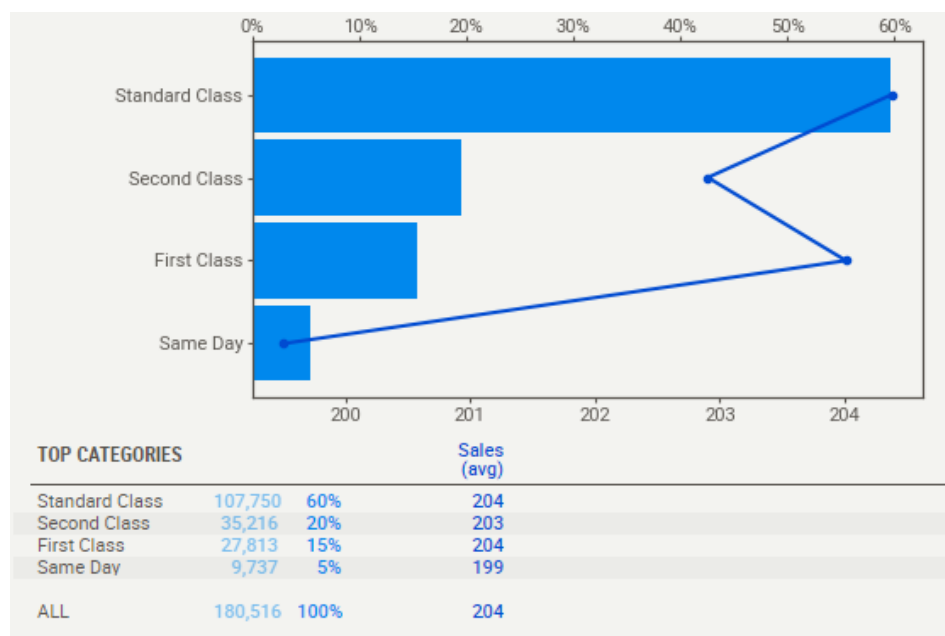
همان‌طور که از نمودار فوق مشاهده می‌شود، فروشگاه‌های ورزشی مخصوص هواداران، پوشاک و گلف حدود ۸۲ درصد دپارتمان‌های فروشگاه را تشکیل می‌دهند سفارش در آن‌ها ثبت می‌شود.



نمودار (۶-۸) نمودار فراوانی بازار محل تحویل سفارش

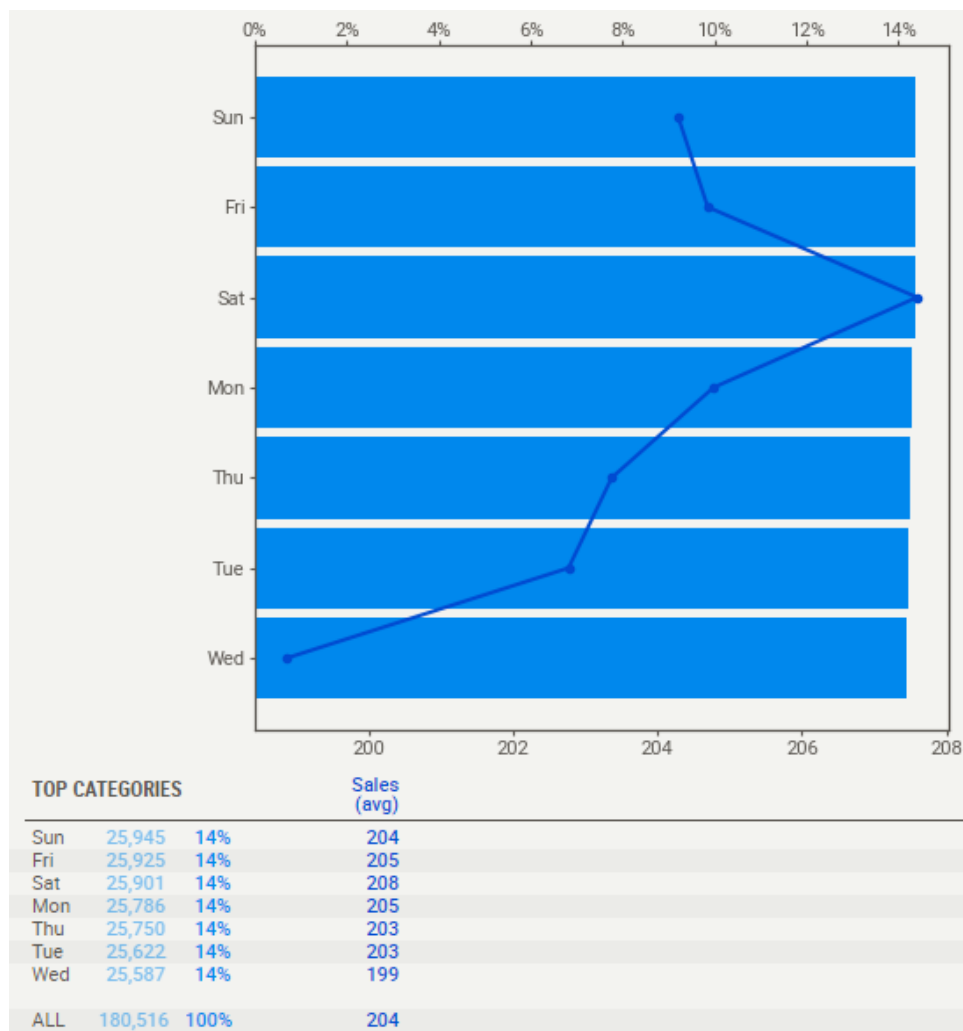
نمودار فوق، فراوانی بازارهای محل تحویل سفارش را نشان می‌دهد. همان‌طور که در این نمودار قابل مشاهده است، بیشترین سفارش‌ها در بازار آمریکای لاتین تحویل داده شده‌اند که ۲۹ درصد سفارش‌ها را به خود اختصاص داده است. ۲۸ درصد سفارش‌ها در بازار اروپا تحویل داده شده‌اند. همچنین ۲۳ درصد سفارش‌ها نیز در بازار آسیا و اقیانوسیه توسط مشتریان تحویل گرفته شده‌اند. ایالات متحده و آفریقا نیز ۲۰ درصد بازار محل تحویل سفارش‌ها را به خود اختصاص داده‌اند. گستردگی بازارها و درصد سفارش‌ها نشان‌دهنده عملکرد وسیع فروشگاه در بازارهای جهانی است. پراکندگی سفارش‌ها خود می‌تواند دلیلی بر ارسال با تأخیر سفارشات باشد که در این مجموعه داده در زمان تحویل تأخیر دارند.

در ادامه بررسی فراوانی حالت‌های حمل‌ونقل سفارش‌ها می‌پردازیم.



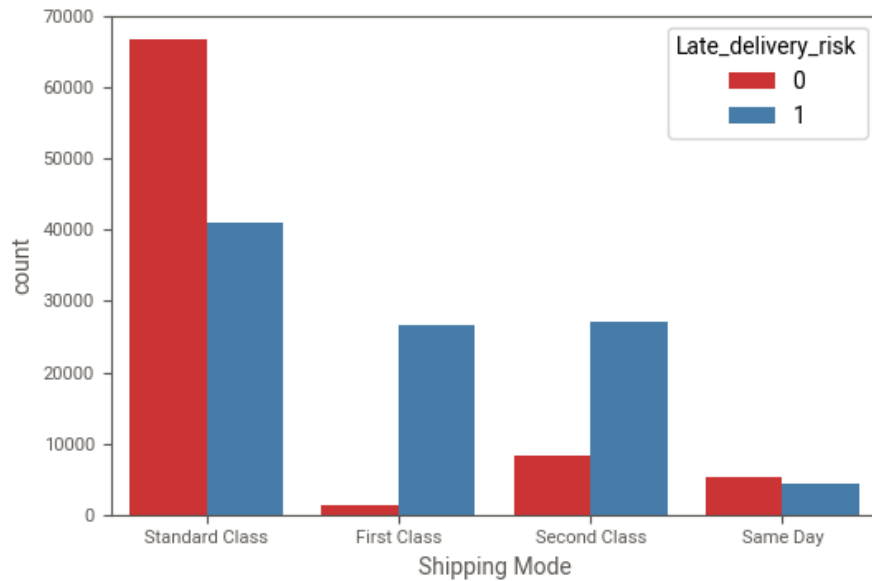
نمودار (۶-۹) نمودار فراوانی حالت‌های حمل‌ونقل سفارش‌ها

باتوجه به نمودار فوق، ۶۰ درصد از سفارش‌ها با کلاس استاندارد حمل‌ونقل آن‌ها صورت گرفته است. ۲۰ درصد سفارش‌ها با حمل‌ونقل درجه دو و ۱۵ درصد سفارش‌ها با حمل‌ونقل درجه یک تحویل داده شده‌اند. تنها ۵ درصد سفارش‌ها در همان روز سفارش تحویل مشتری داده شده‌اند. همچنین باتوجه به نمودار زیر، می‌توان فهمید که توزیع سفارش‌ها در روزهای هفته تقریباً به صورت یکنواخت بوده و روزانه بین ۲۵۹۰۰ الی ۲۶۰۰۰ سفارش در این فروشگاه ثبت می‌شود.



نمودار (۶-۱۰) نمودار توزیع سفارش‌ها فروشگاه طی روزهای هفته

نمودار زیر حالت‌های حمل‌ونقل سفارش‌ها را به تفکیک وجود ریسک ارسال با تأخیر نمایش می‌دهد. همان‌طور که مشاهده می‌شود، سفارش‌ها با حمل‌ونقل درجه یک و درجه دو، ریسک ارسال با تأخیر در آن‌ها بیشتر است. اما در حالت حمل‌ونقل کلاس استاندارد، در بیشتر سفارش‌ها ریسک وجود ارسال با تأخیر را ندارند.



نمودار (۶-۱۱) نمودار حالت‌های حمل‌ونقل سفارش‌ها به تفکیک وجود ریسک ارسال با تأخیر

۶-۲- استخراج قوانین انجمنی

الگوریتم Apriori برای استخراج قوانین انجمنی در این مطالعه استفاده شده است. حداقل حد آستانه پشتیبانی، ۰.۳ در نظر گرفته شده است. جدول زیر تعدادی از قوانین را نشان می‌دهد که بیشترین تکرار را در مجموعه داده داشته‌اند. در واقع بیشترین میزان پشتیبانی را دارند.

جدول (۶-۱) ترتیب قوانین انجمنی استخراج شده بر اساس معیار پشتیبانی از زیاد به کم

support	itemsets
5 0.615696	(Customer Country_USA)
9 0.596900	(Shipping Mode_Standard Class)
10 0.548295	(Delivery Status_Late delivery, Late_delivery_...
0 0.548295	(Late_delivery_risk)
2 0.548295	(Delivery Status_Late delivery)
6 0.517971	(Customer Segment_Consumer)
4 0.384304	(Customer Country_Puerto Rico)
1 0.383861	(Type_DEBIT)
3 0.369884	(Customer City_Caguas)
14 0.369884	(Customer City_Caguas, Customer Country_Puerto...
16 0.366433	(Shipping Mode_Standard Class, Customer Countr...
11 0.337837	(Late_delivery_risk, Customer Country_USA)
13 0.337837	(Delivery Status_Late delivery, Customer Count...
18 0.337837	(Delivery Status_Late delivery, Late_delivery_...
8 0.329550	(Order Status_COMPLETE)
12 0.329550	(Order Status_COMPLETE, Type_DEBIT)
15 0.320348	(Customer Segment_Consumer, Customer Country_USA)
17 0.309950	(Customer Segment_Consumer, Shipping Mode_Stan...
7 0.303508	(Customer Segment_Corporate)

همان طور که مشاهده می شود، کشور ایالات متحده به عنوان کشور محل سفارش مشتری و حالت حمل و نقل استاندارد، در میان موارد تکی بیشترین تعداد را داشته اند. جدول زیر، قوانین انجمنی استخراج شده از جدول فوق را نشان می دهد که معیارهای Confidence و Lift در آن محاسبه شده است. ترتیب این قوانین بر اساس بیشترین میزان معیار Lift می باشد.

جدول (۶-۲) قوانین انجمنی استخراج شده از مجموعه داده، ترتیب بر اساس معیار Lift از زیاد به کم

	antecedents	consequents	antecedent support	consequent support	support	confidence	lift	leverage	conviction	zhangs_metric
4	(Order Status_COMPLETE)	(Type_DEBIT)	0.329550	0.383861	0.329550	1.000000	2.605112	0.203049	inf	0.918993
5	(Type_DEBIT)	(Order Status_COMPLETE)	0.383861	0.329550	0.329550	0.858514	2.605112	0.203049	4.738628	1.000000
8	(Customer City_Caguas)	(Customer Country_Puerto Rico)	0.369884	0.384304	0.369884	1.000000	2.602107	0.227736	inf	0.977116
9	(Customer Country_Puerto Rico)	(Customer City_Caguas)	0.384304	0.369884	0.369884	0.962478	2.602107	0.227736	16.793327	1.000000
0	(Delivery Status_Late delivery)	(Late_delivery_risk)	0.548295	0.548295	0.548295	1.000000	1.823836	0.247668	inf	1.000000
20	(Late_delivery_risk)	(Delivery Status_Late delivery, Customer Count...	0.548295	0.337837	0.337837	0.616159	1.823836	0.152603	1.725099	1.000000
19	(Delivery Status_Late delivery)	(Late_delivery_risk, Customer Country_USA)	0.548295	0.337837	0.337837	0.616159	1.823836	0.152603	1.725099	1.000000
18	(Late_delivery_risk, Customer Country_USA)	(Delivery Status_Late delivery)	0.337837	0.548295	0.337837	1.000000	1.823836	0.152603	inf	0.682166
17	(Delivery Status_Late delivery, Customer Count...	(Late_delivery_risk)	0.337837	0.548295	0.337837	1.000000	1.823836	0.152603	inf	0.682166
1	(Late_delivery_risk)	(Delivery Status_Late delivery)	0.548295	0.548295	0.548295	1.000000	1.823836	0.247668	inf	1.000000
11	(Customer Segment_Consumer)	(Customer Segment_Consumer)	0.615696	0.517971	0.320348	0.520303	1.004502	0.001436	1.004861	0.011662
10	(Customer Segment_Consumer)	(Customer Country_USA)	0.517971	0.615696	0.320348	0.618468	1.004502	0.001436	1.007265	0.009298
14	(Customer Segment_Consumer)	(Shipping Mode_Standard Class)	0.517971	0.596900	0.309950	0.598394	1.002502	0.000774	1.003719	0.005178
15	(Shipping Mode_Standard Class)	(Customer Segment_Consumer)	0.596900	0.517971	0.309950	0.519267	1.002502	0.000774	1.002696	0.006192
16	(Delivery Status_Late delivery, Late_delivery_...	(Customer Country_USA)	0.548295	0.615696	0.337837	0.616159	1.000753	0.000254	1.001207	0.001665
7	(Customer Country_USA)	(Delivery Status_Late delivery)	0.615696	0.548295	0.337837	0.548708	1.000753	0.000254	1.000914	0.001957
6	(Delivery Status_Late delivery)	(Customer Country_USA)	0.548295	0.615696	0.337837	0.616159	1.000753	0.000254	1.001207	0.001665
3	(Customer Country_USA)	(Late_delivery_risk)	0.615696	0.548295	0.337837	0.548708	1.000753	0.000254	1.000914	0.001957
2	(Late_delivery_risk)	(Customer Country_USA)	0.548295	0.615696	0.337837	0.616159	1.000753	0.000254	1.001207	0.001665

در سطر اول جدول مشاهده می گردد، وضعیت سفارش کامل با نوع معامله و کالتی به هم وابسته هستند و در کنار یکدیگر می آیند. سایر موارد و قوانین استخراج شده نیز در فایل پیوست کدها قابل مشاهده است.

۶-۳- ارزیابی و مقایسه عملکرد الگوریتم های یادگیری ماشین

در جدول زیر، مقایسه ارزیابی الگوریتم های اجرا شده روی داده های دیده نشده توسط مدل صورت گرفته است. ابتدا مقایسه الگوریتم های رگرسیون برای پیش بینی میزان فروش هر سفارش را بررسی می کنیم.

جدول (۶-۳) مقایسه الگوریتم های رگرسیون یادگیری ماشین برای پیش بینی فروش هر سفارش

	MLA used	Train Accuracy	Test Accuracy	MAE	MSE	RMSE	R ² Score	Adjusted R ² Score
0	LinearRegression	1.0000	1.0000	0.000593	0.000002	0.001563	1.000000	1.000000
1	Ridge	1.0000	1.0000	0.000593	0.000002	0.001563	1.000000	1.000000
2	Lasso	1.0000	1.0000	0.012694	0.000407	0.020177	1.000000	1.000000
3	DecisionTreeRegressor	1.0000	1.0000	0.012253	0.656473	0.810230	0.999963	0.999963
5	BayesianRidge	1.0000	1.0000	0.000593	0.000002	0.001563	1.000000	1.000000
4	RandomForestRegressor	0.8585	0.8636	38.873632	2444.683456	49.443740	0.863609	0.862842

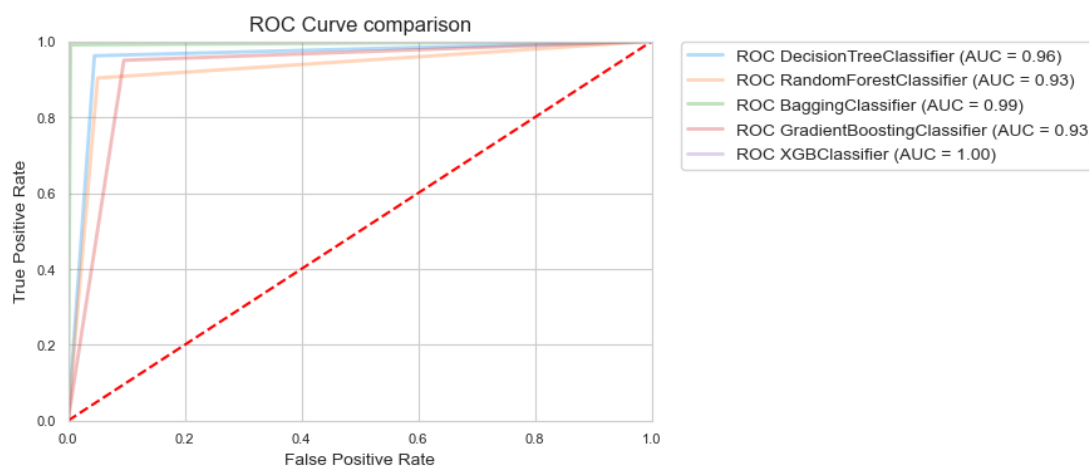
همان طور که مشاهده می شود، از ۶ مدل پیاده سازی شده، تنها الگوریتم رگرسیون جنگل تصادفی عملکرد ضعیف تری داشته است. الگوریتم های رگرسیون خطی، رگرسیون ستیغی و رگرسیون بیزی با معیار مربع R تنظیم شده یکسان ۱، بهترین مدل ها برای پیش بینی فروش هر سفارش هستند. در ادامه، به مقایسه الگوریتم های دسته بندی برای پیش بینی وجود ریسک ارسال با تأخیر برای هر سفارش می پردازیم.

جدول (۴-۶) مقایسه الگوریتم‌های دسته‌بندی یادگیری ماشین برای پیش‌بینی وجود ریسک ارسال با تأخیر

	MLA used	Train Accuracy	Test Accuracy	Precision	Recall	AUC	F1-Score
4	XGBClassifier	1.0000	0.9998	0.999695	1.000000	0.999818	0.999847
2	BaggingClassifier	0.9998	0.9838	0.993476	0.976635	0.984490	0.984983
0	DecisionTreeClassifier	1.0000	0.9587	0.962169	0.961924	0.958391	0.962047
3	GradientBoostingClassifier	0.9317	0.9292	0.922267	0.950013	0.927221	0.935934
1	RandomForestClassifier	1.0000	0.9244	0.955167	0.903385	0.926387	0.928555

مشاهده می‌شود که الگوریتم XGboost با امتیاز $F1$ نزدیک به ۱ به‌خوبی پیش‌بینی را انجام داده است و سایر الگوریتم‌ها نیز در محدوده ۰.۹۲ تا ۰.۹۹ قرار دارند که این موضوع نشان از کارایی بسیار خوب الگوریتم‌های یادگیری ماشین در این زمینه دارد.

همچنین نمودار زیر، منحنی ROC الگوریتم‌های دسته‌بندی بکار گرفته شده را نشان می‌دهد. مساحت زیر منحنی مربوط به مدل XGboost بیشترین مقدار را در بین الگوریتم‌های پیاده‌سازی شده در اختیار دارد.



نمودار (۴-۶) منحنی ROC برای الگوریتم‌های دسته‌بندی بکار گرفته شده

۷- نتیجه‌گیری و پیشنهاد

در این مطالعه، برای تجزیه و تحلیل زنجیره تامین یک فروشگاه خرده‌فروشی با فروش حضوری و اینترنتی و پیاده‌سازی الگوریتم‌های یادگیری ماشین، ما از مجموعه داده اطلاعات زنجیره تامین یک فروشگاه اینترنتی که توسط شرکت دیتاکو^{۱۰۳} بین سال‌های ۲۰۱۵ تا ۲۰۱۸، منتشر شده است، استفاده کرده‌ایم. تجزیه و تحلیل اکتشافی داده‌ها و نتایج این مطالعه نشان می‌دهد:

- ویژگی‌های مجموع فروش انجام شده به ازای هر مشتری، مقدار تخفیف کالای سفارشی، قیمت محصولات بدون تخفیف، مبلغ کل به ازای هر سفارش و قیمت محصول با میزان فروش همبستگی بالایی (بین ۰.۷۵ تا ۱) دارند. در واقع با افزایش هر کدام از این ویژگی‌ها می‌توان انتظار داشت که میزان فروش افزایش یابد. یکی از راه‌های افزایش فروش می‌تواند ایجاد کمپین‌های تبلیغاتی دوره‌ای و افزایش میزان تخفیفات کالاها باشد.
- در وضعیت تحویل سفارش‌ها باید به طور ویژه به آن پرداخته شود و با بررسی علل ریشه‌ای هر کدام از سفارش‌ها، پروژه‌های بهبود را برای حمل و نقل سفارش‌ها با دقت بیشتری دنبال نمود.
- باتوجه به تعداد سفارش‌ها انجام شده می‌توان متوجه شده که محصولات کفش ورزشی، کفش مردانه و پوشاک زنانه نزدیک به ۴۰ درصد سفارش‌ها فروشگاه را شامل می‌شوند که فروشگاه باتوجه به این موضوع باید برای فراهم کردن این محصولات برنامه‌ریزی دقیق‌تری بر اساس سلیقه مصرف‌کننده داشته باشد.
- الگوریتم‌های یادگیری ماشین رگرسیون و دسته‌بندی به ترتیب عملکرد فوق‌العاده‌ای در پیش‌بینی میزان فروش هر سفارش و وجود ریسک ارسال با تأخیر سفارش داشته‌اند. این موضوع نشان می‌دهد که استفاده از الگوریتم‌های یادگیری ماشین می‌تواند عملکرد بسیار مؤثری در حوزه زنجیره تامین فروشگاه اینترنتی داشته باشند.
- در مطالعات آینده، پیشنهاد می‌شود از داده‌های حجیم‌تر برای پیاده‌سازی الگوریتم‌های یادگیری ماشین استفاده شود. همچنین به دلیل بهبود پیوسته و مداوم الگوریتم‌ها، می‌توان از الگوریتم‌های موجود دیگر مانند Adaboost برای پیش‌بینی استفاده کرد و نتایج را با نتایج این پژوهش مقایسه نمود.
- حوزه زنجیره تامین، در حال حاضر به عنوان یک ترند جهانی به طور پیوسته در حال بهبود است و به کارگیری هوش مصنوعی در فرایندهای مختلف زنجیره تامین می‌تواند بسیار مؤثر واقع گردد.

^{۱۰۳} Data-Co

استفاده از حوزه‌های یادگیری عمیق^{۱۰۴}، پردازش زبان طبیعی^{۱۰۵} و بینایی ماشین^{۱۰۶} در مطالعات آینده می‌تواند مورد توجه قرار گیرد.

^{۱۰۴} Deep Learning

^{۱۰۵} Natural Language Process

^{۱۰۶} Computer Vision

١. Li, L., et al., *Data-driven online service supply chain: a demand-side and supply-side perspective*. Journal of Enterprise Information Management, ٢٠٢١. ٣٤(١): p. ٣٦٥-٣٨١.
٢. Pereira, M.M. and E.M. Frazzon, *A data-driven approach to adaptive synchronization of demand and supply in omni-channel retail supply chains*. International Journal of Information Management, ٢٠٢١. ٥٧: p. ١٠٢١٦٥.
٣. Rodrigues, V.S., et al., *Measurement, mitigation and prevention of food waste in supply chains: An online shopping perspective*. Industrial Marketing Management, ٢٠٢١. ٩٣: p. ٥٤٥-٥٦٢.
٤. Krishna, A., et al. *Sales-forecasting of retail stores using machine learning techniques*. in *2018 3rd international conference on computational systems and information technology for sustainable solutions (CSITSS)*. ٢٠١٨. IEEE.
٥. Beheshti-Kashi, S., et al., *A survey on retail sales forecasting and prediction in fashion markets*. Systems Science & Control Engineering, ٢٠١٥. ٣(١): p. ١٥٤-١٦١.
٦. Oke, A. and M. Gopalakrishnan, *Managing disruptions in supply chains: A case study of a retail supply chain*. International journal of production economics, ٢٠٠٩. ١١٨(١): p. ١٦٨-١٧٤.
٧. Jiang, H., J. Ruan, and J. Sun. *Application of machine learning model and hybrid model in retail sales forecast*. in *٢٠٢١ IEEE 6th international conference on big data analytics (ICBDA)*. ٢٠٢١. IEEE.
٨. Yao, J. and M. Gu, *Optimization analysis of supply chain resource allocation in customized online shopping service mode*. Mathematical Problems in Engineering, ٢٠١٥. ٢٠١٥.
٩. Yao, J., *Supply chain resources integration optimisation in B2C online shopping*. International Journal of Production Research, ٢٠١٧. ٥٥(١٧): p. ٥٠٧٩-٥٠٩٤.
١٠. Kück, M. and M. Freitag, *Forecasting of customer demands for production planning by local k-nearest neighbor models*. International Journal of Production Economics, ٢٠٢١. ٢٣١: p. ١٠٧٨٣٧.
١١. Shilong, Z. *Machine learning model for sales forecasting by using XGBoost*. in *2021 IEEE International Conference on Consumer Electronics and Computer Engineering (ICCECE)*. ٢٠٢١. IEEE.
١٢. Sajawal, M., et al., *A Predictive Analysis of Retail Sales Forecasting using Machine Learning Techniques*. Lahore Garrison University Research Journal of Computer Science and Information Technology, ٢٠٢٢. ٦(٠٤): p. ٣٣-٤٥.
١٣. Lee, I. and G. Mangalaraj, *Big data analytics in supply chain management: a systematic literature review and research directions*. Big Data and Cognitive Computing, ٢٠٢٢. ٦(١): p. ١٧.
١٤. Li, J., A. Ghadge, and M.K. Tiwari, *Impact of replenishment strategies on supply chain performance under e-shopping scenario*. Computers & Industrial Engineering, ٢٠١٦. ١٠٢: p. ٧٨-٨٧.

- .۱۵ Chen, I.-F. and C.-J. Lu, *Sales forecasting by combining clustering and machine-learning techniques for computer retailing*. Neural Computing and Applications, ۲۰۱۷. ۲۸: p. ۲۶۳۳-۲۶۴۷.
- .۱۶ Catal, C., et al., *Benchmarking of regression algorithms and time series analysis techniques for sales forecasting*. Balkan Journal of Electrical and Computer Engineering, ۲۰۱۹. ۷(۱): p. ۲۰-۲۶.
- .۱۷ Arif ,M.A.I., et al. *Comparison study: Product demand forecasting with machine learning for shop*. in *2019 8th International Conference System Modeling and Advancement in Research Trends (SMART)*. ۲۰۱۹. IEEE.
- .۱۸ Oosthuizen, K., et al., *Artificial intelligence in retail: The AI-enabled value chain*. Australasian Marketing Journal, ۲۰۲۱. ۲۹(۳): p. ۲۶۴-۲۷۳.
- .۱۹ Wu, Y., et al., *Low-carbon decision-making model of online shopping supply chain considering the O2O model*. Journal of Retailing and Consumer Services, ۲۰۲۱ . :۵۹p. ۱۰۲۳۸۸.
- .۲۰ Yu, Y., et al., *E-commerce logistics in supply chain management: Practice perspective*. Procedia Cirp, ۲۰۱۶. ۵۲: p. ۱۷۹-۱۸۵.
- .۲۱ Abdirad, M. and K. Krishnan, *Examining the impact of E-supply chain on service quality and customer satisfaction: a case study*. International Journal of Quality and Service Sciences, ۲۰۲۲. ۱۴(۲): p. ۲۷۴-۲۹۰.
- .۲۲ Chen, F. and T. Ou, *Sales forecasting system based on Gray extreme learning machine with Taguchi method in retail industry*. Expert Systems with Applications, ۲۰۲۱. ۲۸(۳): p. ۱۳۳۶-۱۳۴۵.
- .۲۳ Liu, Y., et al. *An Aggregate Store Sales Forecasting Framework based on ConvLSTM*. in *2021 The 5th International Conference on Compute and Data Analysis*. ۲۰۲۱.
- .۲۴ Andrejić, M., *Modeling Retail Supply Chain Efficiency :Exploration and Comparative Analysis of Different Approaches*. Mathematics, ۲۰۲۳. ۱۱(۷): p. ۱۵۷۱.