

مقدمه ای بر یادگیری ماشین

سولویگ بدیلوا^{۱،*}، بلاژ بنفای^۱، فابیان بیرزله^۱، یاکوف آی. داویدوف^۱، لوسی هاجینسون^۱،
تونی کام-تونگ^۱، جولیان سیبور-پولسترا^۱، برنهارد استیرت^۱، وجیتایو دیوید ژانگ^۱

در چند سال گذشته، یادگیری ماشین (ML) و هوش مصنوعی موج جدیدی از تبلیغات را به خود دیده اند که ناشی از حجم عظیم و روزافزون داده ها و قدرت محاسباتی و همچنین کشف الگوریتم های یادگیری بهبود یافته است. با این حال، ایده یادگیری مفاهیم انتزاعی از داده ها توسط کامپیوتر و اعمال آنها در موقعیت های ناشناخته، چیز جدیدی نیست و حداقل از دهه ۱۹۵۰ وجود داشته است. بسیاری از این اصول اساسی برای جامعه فارماکومتریک و فارماکولوژی بالینی بسیار آشنا هستند. در این مقاله، می خواهیم ایده های اساسی یادگیری ماشین را به این جامعه معرفی کنیم تا خوانندگان ابزارهای ضروری مورد نیاز برای درک انتشارات مربوط به این موضوع را به دست آورند. اگرچه ما به جزئیات و پیشینه نظری نخواهیم پرداخت، اما هدف ما این است که خوانندگان را به ادبیات مربوطه ارجاع دهیم و کاربردهای یادگیری ماشین را در زیست شناسی مولکولی و همچنین زمینه های فارماکومتریک و فارماکولوژی بالینی در چشم انداز قرار دهیم.

تعاملات این پیچیدگی می تواند چالش هایی را برای تفسیر مدل ایجاد کند (به اصطلاح مشکل "جعبه سیاه"). رویکردهایی که معمولاً در کاربردهای فارماکومتری استفاده می شوند، در فرهنگ قرار می گیرند، که در آن یک مدل اساسی بر اساس اصول دارویی و درک خواص دارو فرض می شود. چنین مدل هایی معمولاً از نظر فیزیولوژیکی قابل تفسیر هستند. اکثر رویکردهای یادگیری ماشین (ML) در فرهنگ قرار می گیرند، که در آن هیچ مدل صریحی مشخص نشده است و یک کامپیوتر مسئول شناسایی ارتباطات در داده های مشاهده شده است. تفسیر این مدل ها از نظر فیزیولوژیکی دشوار است، با این حال، در طول سال های پیشرفت قابل توجهی در تفسیرپذیری مدل های ML حاصل شده است.^{۱،۲} امروزه، بسیاری از جنبه های یک مدل جعبه سیاه را می توان با استفاده از ابزارهای مناسب تفسیر کرد.^۳ در این مقاله، هدف ما پشتیبانی از خوانندگان برای توسعه شهود مورد نیاز برای درک چگونگی یادگیری کامپیوترها یا کمک به انسان ها در شناسایی الگوها در داده ها است. ایده های بنیادی یادگیری ماشین برجسته شده اند، اما جزئیات و پیشینه نظری روش های یادگیری ماشین موجود را شرح نمی دهیم. خوانندگان علاقه مند را به مقالات یا کتاب های دیگر، مانند «عناصر یادگیری آماری» ارجاع می دهیم.^۴ (که به عنوان ESL شناخته می شود) و ما به نمونه هایی از کاربرد آنها در زیست شناسی مولکولی، کشف دارو، توسعه دارو و داروشناسی بالینی اشاره می کنیم.

ما ابتدا مفاهیم نقاط داده، ویژگی ها، فضاها و ویژگی و معیارهای شباهت را معرفی می کنیم و سپس عمیق تر به دو حوزه اصلی یادگیری ماشین، یعنی یادگیری بدون نظارت و یادگیری تحت نظارت، می پردازیم و جنبه های کلیدی و مثال ها را بررسی می کنیم. در مورد یادگیری بدون نظارت، کامپیوترها وظیفه دارند الگوهای ناشناخته در داده ها را بدون دانش از پیش موجود مانند گروه های کلاس ها شناسایی کنند، در حالی که در مورد یادگیری تحت نظارت، کامپیوترها وظیفه دارند یاد بگیرند که چگونه کلاس یا مقدار نقاط داده هنوز مشاهده نشده را بر اساس یک مفهوم (که اغلب "مدل" نیز نامیده می شود) که از یک مجموعه داده آموزشی استخراج شده است، پیش بینی کنند. شکل ۱ طبقه بندی روش های مختلف شرح داده شده در این مقاله را نشان می دهد و می تواند

ظهور دسترسی به داده ها و رشد قدرت محاسباتی، همراه با ورود روش های نوین یادگیری، منجر به پیشرفت های چشمگیری در بسیاری از حوزه های علمی شده است. این شامل تحقیقات بیولوژیکی و بالینی می شود که کاربردهای آن از زیست شناسی مولکولی گرفته تا ... برای تجزیه و تحلیل داده های تصویری^۵ و عمل بالینی^۶ با این حال، ایده یادگیری مفاهیم انتزاعی توسط کامپیوتر - مانند کاری که انسان ها دائماً انجام می دهند - حداقل از دهه ۱۹۵۰، زمانی که اولین شبکه های عصبی ساخته شدند، وجود داشته است.^۷ توسعه داده شدند. حتی قبل از آن، روش های دیگری مانند آمار بیزی و زنجیره های مارکوف با ایده ای مشابه مورد استفاده قرار می گرفتند. بسیاری از این روش ها با قراردادهای نامگذاری مختلف برای جامعه فارماکومتریک و فارماکولوژی بالینی شناخته شده اند. در سمت چپ، اصطلاحات یادگیری ماشین و در سمت راست، نامگذاری معمول آمار (بر اساس <https://statweb.stanford.edu/~tibs/stat315a/glossary.pdf>) را نشان می دهیم:

- شبکه، نمودارها ⇒ مدل
- وزنه ها ⇒ پارامترها
- یادگیری ⇒ مناسب سازی
- تعمیم ⇒ عملکرد مجموعه آزمون
- یادگیری تحت نظارت ⇒ رگرسیون یا طبقه بندی
- یادگیری بدون نظارت ⇒ تخمین چگالی، خوشه بندی
- ویژگی ها ⇒ متغیرهای کمکی یا توضیحی

تفاوت اصلی با رویکردهای سنتی تر، تا حد زیادی در دو فرهنگ متمایز مدل سازی آماری نهفته است. بریمن تقریباً دو دهه پیش از این موضوع طفره رفته بود. در اینجا، ما تعریف او را با گنجاندن مدل های فیزیولوژیکی در یکی از فرهنگ ها گسترش می دهیم. به طور خاص، فرهنگ ۱ شامل تعیین مدلی برای توصیف داده های مشاهده شده است و فرهنگ ۲ با اتخاذ یک رویکرد مدل سازی الگوریتمی، قصد دارد مسئله را حل کند، بنابراین ذاتاً منجر به مدل هایی با تعداد بیشتری از پارامترهای آزاد و پیچیدگی های بیشتر می شود.

علوم دارویی، تحقیقات و توسعه اولیه داروسازی روشه (pRED)، مرکز نوآوری روشه بازل، بازل، سوئیس.
* مکاتبه: Solveig Badillo (solveig.badillo@roche.com) نویسندگان به ترتیب حروف الفبا. همه نویسندگان به طور مساوی مشارکت داشتند.
[اصلاحیه در تاریخ ۶ مارس ۲۰۲۰، پس از اولین انتشار آنلاین اضافه شد: متن مشارکت نویسنده اضافه شد].
آ.س.بی. در زمان نگارش این دست نوشته، در استخدام گروه سولادیس بود.
دریافت شده در ۸ اکتبر ۲۰۱۹؛ پذیرفته شده در ۱۵ ژانویه ۲۰۲۰. doi:10.1002/cpt.1796



داده ها و ویژگی ها

در یادگیری ماشین، ما با داده ها و مجموعه داده ها سروکار داریم. یک مجموعه داده از چندین نقطه داده (که گاهی، نمونه نیز نامیده می شوند) تشکیل شده است، که در آن هر

میانگین ویژگی در تمام نمونه هایی که در آنها تعریف شده است. با این حال، این گاهی اوقات می تواند باعث بیش برآزش شود (همچنین به بخش «معیارهای عملکرد و مسئله بیش برآزش» مراجعه کنید). همچنین بررسی دقیق هرگونه سوگیری در داده ها (مثلاً سوگیری انتخاب) ضروری است. ترجیحاً، نمونه های ML باید یک زیرمجموعه تصادفی بدون سوگیری از جمعیت باشند. در عمل، این امر به ندرت اتفاق می افتد و برخی سوگیری ها در داده ها وجود دارد. این سوگیری های می توانند بر توانایی مدل برای تعمیم فراتر از مجموعه داده های آموزشی (و حتی مجموعه داده های آزمایشی اگر هر دو سوگیری مشابهی داشته باشند) تأثیر بگذارند. نمونه ای از چنین مسئله تعمیمی، مدلی است که قرار است یاد بگیرد چگونه گرگ را از هاسکی بر اساس ویژگی های حیوان تشخیص دهد، اما در نهایت معلوم می شود که به سادگی تکه های برف را روی عکس تشخیص می دهد. رویکردهای مختلفی برای کاهش سوگیری وجود دارد (مثلاً می توان نمونه ها یا ویژگی های سوگیرانه را کم وزن کرد یا به طور کامل حذف کرد).^{۱۲} به طور خاص، نمرات گرایش هنگام تخمین اثر یک مداخله درمانی مفید هستند.^{۱۳} بررسی اهمیت ویژگی، اطلاعات ارزشمندی در مورد بزرگی و تأثیر بایاس ارائه می دهد.^{۶،۷} که توصیه می شود برای بررسی قابلیت اعتماد مدل های یادگیری ماشینی استفاده شود.

بسیاری از مجموعه داده های طبقه بندی بالینی نامتوازن هستند، به این معنی که یک یا چند کلاس کمتر از حد واقعی نمایش داده شده اند. این می تواند برای بسیاری از الگوریتم های یادگیری ماشینی، از جمله شبکه های عصبی مصنوعی و روش های تقویت گرادیان، مشکلاتی ایجاد کند. یک راه برای کاهش این مشکل، نمونه گیری کمتر/بیش از حد به ترتیب از کلاس اکثریت/اقلیت، یا تنظیم هزینه طبقه بندی نادرست در تابع هدف است.^{۱۴} در نهایت، برای بسیاری از کاربردها، تعریف یک معیار شباهت یا فاصله بین دو نقطه داده در فضای ویژگی مهم است. ساده ترین معیار فاصله، فاصله اقلیدسی است:

$$d = \sqrt{\sum_{i=1}^n (x_i - y_i)^2} \quad (1)$$

(الف، ب) = $\sqrt{\sum_{i=1}^n (x_i - y_i)^2}$ (یک من ب من) = n

بین بردارهای ویژگی عددی دو نقطه داده/الف و ب، برای ویژگی ها $1 \dots n$ ، اما بسته به نوع داده ای که با آن سروکار داریم، می تواند معیارهای فاصله یا شباهت بسیار پیچیده تری مانند شباهت کسینوسی وجود داشته باشد.^{۱۵} یا امتیازهای شباهت دو توالی بیولوژیکی.^{۱۶}

غذاهای اصلی

- تبدیل داده های ورودی و مهندسی ویژگی ها ممکن است مدل را بهبود بخشد.
- داده های گمشده نیاز به جانشینی دارند.
- سوگیری های موجود در داده ها باید به دقت بررسی شوند.
- مجموعه داده های نامتوازن نیاز به اصلاح مدل دارند.
- معیارهای معنادار شباهت بین نمونه ها باید تعریف شوند.

یادگیری بدون نظارت

در تحلیل اکتشافی داده ها، ما اغلب «برچسب های» واقعی را نمی دانیم، یا ممکن است بخواهیم الگوهای طبیعی در حال ظهور در داده ها را بررسی کنیم. برای این منظور، می توانیم از روش های یادگیری بدون نظارت، مانند خوشه بندی، تشخیص الگوهای مکرر و کاهش ابعاد استفاده کنیم. در اینجا، ما به طور خاص تمرکز خواهیم کرد.

امتیاز). این ویژگی ها می توانند دسته بندی شده (مقادیر از پیش تعریف شده بدون ترتیب خاص مانند مرد و زن)، ترتیبی (مقادیر از پیش تعریف شده که دارای ترتیب ذاتی مانند مرحله بیماری هستند) یا عددی (مثلاً مقادیر واقعی) باشند. برای یک بیمار در یک محیط بالینی، این ویژگی ها می توانند (ترکیبی از) اطلاعات جمعیت شناختی بیمار، سابقه بیماری، نتایج آزمایش خون یا معیارهای پیچیده تر و با ابعاد بالا، مانند پروفایل های بیان ژن در یک بافت خاص یا تمام پلی مورفیسم های تک نوکلئوتیدی که نشان دهنده ژنوم منحصر به فرد بیمار هستند، باشند. هر ویژگی، یک بُعد از فضای ویژگی را نشان می دهد و مقدار مشخص یک ویژگی برای یک نقطه داده خاص، آن نقطه را در مکانی تعریف شده در این بُعد از فضا قرار می دهد. در مجموع، تمام مقادیر تمام ویژگی های یک نقطه داده، بردار ویژگی نامیده می شود. هرچه ویژگی های بیشتری برای مجموعه داده ها جمع آوری کرده باشیم، بردار ویژگی حاصل و فضای ویژگی، ابعاد بالاتری خواهند داشت. بدیهی است که با افزایش ابعاد، تجسم تمام ابعاد فضای ویژگی، دشواری شود و ما باید برای شناسایی الگوهای مربوطه به کامپیوتر تکیه کنیم یا از روش های کاهش ابعاد استفاده کنیم، همانطور که بعداً در بخش «کاهش ابعاد» توضیح داده خواهد شد.

دارو شناسان بالینی معمولاً با داده های طولی، مانند پروفایل های فارماکوکینتیک (PK) و فارماکودینامیک (PD)، که در آن ها وابستگی به زمان نقش محوری دارد، آشنا هستند. در واقع، مدل های مورد استفاده در فارماکومتری مبتنی بر معادلاتی هستند که می توانند بر اساس فیزیولوژی و فارماکولوژی توجیه شوند و بینش هایی در مورد تکامل زمانی سیستم ارائه دهند. این امر، به عنوان مثال، مشابه مسائل فیزیکی مانند پیش بینی آب و هوا است که در آن جریان هوا و دما منجر به رفتار زمانی خاصی از سیستم می شوند. در یادگیری ماشینی، گنجاندن زمان به عنوان یک متغیر پیوسته متمایز در الگوریتم های مربوطه، همچنان چالش برانگیز است و حوزه ای از تحقیقات فعال است. در حال حاضر، چندین گزینه برای گنجاندن داده های وابسته به زمان در مجموعه داده های یادگیری ماشینی وجود دارد: یا به طور مستقیم که در آن هر نقطه زمانی نشان دهنده یک ویژگی است، یا از طریق تبدیلاتی مانند تبدیل فوریه یا B-splines که منجر به ضرایب توابع اساسی می شود که می توانند به عنوان ویژگی در نظر گرفته شوند. از طرف دیگر، شبکه های عصبی بازگشتی (RNN) می توانند برای مدیریت داده های طولی، همانطور که در بخش «شبکه عصبی بازگشتی» ذکر شده است، استفاده شوند. با این حال، همه این رویکردها محدودیت گسسته سازی بُعد زمانی - چه مستقیم و چه غیرمستقیم - را دارند.

اکثر الگوریتم های یادگیری ماشینی برای مدیریت مجموعه داده های با ابعاد بالا طراحی شده اند. از این رو، ویژگی های مشتق شده از داده های موجود اغلب شامل می شوند، مانند داده های تبدیل شده با لگاریتم، حاصلضرب ها و نسبت های ویژگی ها یا ترکیبات پیشرفته تر. چنین تبدیل داده ای یک مرحله پیش پردازش مهم است که می تواند تأثیر عمیقی بر عملکرد مدل داشته باشد. بنابراین، همیشه ایده خوبی است که از دانش و تخصص موجود در حوزه برای دستیابی به ویژگی های مرتبط استفاده شود، فرآیندی که گاهی اوقات به عنوان مهندسی ویژگی شناخته می شود.

کیفیت داده ها نقش حیاتی در یادگیری ماشینی ایفا می کند. روش های یادگیری ماشینی که با دقت انتخاب شده اند و بازرسی بصری، از مقادیر شدید داده های پرت جلوگیری می کنند. با این حال، داده های گمشده می توانند چالش برانگیز باشند. همه روش ها از داده های گمشده پشتیبانی نمی کنند و در چنین مواردی، تبدیل داده ها می تواند به عنوان یک مرحله پیش پردازش مورد نیاز باشد. روش های مختلفی برای جایگذاری داده های گمشده وجود دارد که عملکرد آنها به مجموعه داده ها و روش مورد استفاده بستگی دارد. بدیهی ترین رویکرد برای جانشینی، جایگزینی یک مقدار گمشده با

درمورد خوشه بندی و کاهش ابعاد، زیرا کاربردهای زیادی در زیست شناسی مولکولی و عمل بالینی دارند.

خوشه بندی

هدف از به کارگیری روش های خوشه بندی، شناسایی زیرگروه های مرتبط در یک مجموعه داده معین بدون داشتن فرضیه ای از پیش تعریف شده در مورد ویژگی هایی است که زیرگروه ها ممکن است داشته باشند. به عنوان مثال، در گروهی از بیماران مبتلا به یک بیماری خاص، ممکن است بخواهیم زیرگروه هایی را شناسایی کنیم که نشان دهنده مکانیسم های بیولوژیکی متمایزی هستند که بیماری را بر اساس اقدامات مولکولی انجام شده هدایت می کنند.^{۱۷}

یک خوشه زیرمجموعه ای از داده ها است که به یکدیگر «شبهه» هستند، در حالی که نقاط متعلق به خوشه های مختلف «متفاوت تر» هستند. رویکردهای متعددی برای خوشه بندی وجود دارد که از الگوریتم های اساسی مختلفی برای گروه بندی نقاط داده بر اساس «شباهت» آنها استفاده می کنند. همه آنها مزایا و معایبی دارند و بسته به کاربرد و ویژگی های داده ها، باید با دقت انتخاب شوند.

یک رویکرد ساده برای خوشه بندی یک به معنی خوشه بندی است.^{۱۸} اینجا، تعداد خوشه هایی که باید شناسایی شوند، توسط یک پارامتر انتخابی کاربر از پیش تعریف شده است. هر خوشه توسط یک مرکز خوشه نمایش داده می شود، که یک نقطه داده مصنوعی است که نشان دهنده میانگین (یا میانه) مقدار تمام نقاط اختصاص داده شده به این خوشه است. در ابتدا، یک مراکز خوشه، که به عنوان "دانه ها" شناخته می شوند، به طور تصادفی در فضای ویژگی قرار می گیرند. سپس الگوریتم دو مرحله را تکرار می کند. در مرحله اول ("انتساب")، نقاط داده به خوشه ای که توسط نزدیکترین مرکز نشان داده شده است، اختصاص داده می شوند. در مرحله دوم ("تغییر مرکز")، موقعیت هر مرکز خوشه بر اساس ترکیب خوشه ها پس از مرحله اول به روز می شود. پس از تعدادی تکرار، این معمولاً به یک بهینه محلی همگرا می شود که در آن انتساب خوشه ها تغییر نمی کند یا فقط به میزان کمی تغییر می کند. نتیجه چنین فرآیندی در شکل نشان داده شده است. **شکل ۲** اگرچه این روش شهودی است، اما اشکال اصلی آن این است که معمولاً خوشه بندی به شدت تحت تأثیر مقدار ... قرار می گیرد. کو اغلب تعداد واقعی خوشه ها در داده ها ناشناخته است به طور پیشینی از آنجا که در خوشه بندی به ندرت یک پاسخ درست یا غلط قطعی وجود دارد، بررسی بیشتر خوشه بندی برای

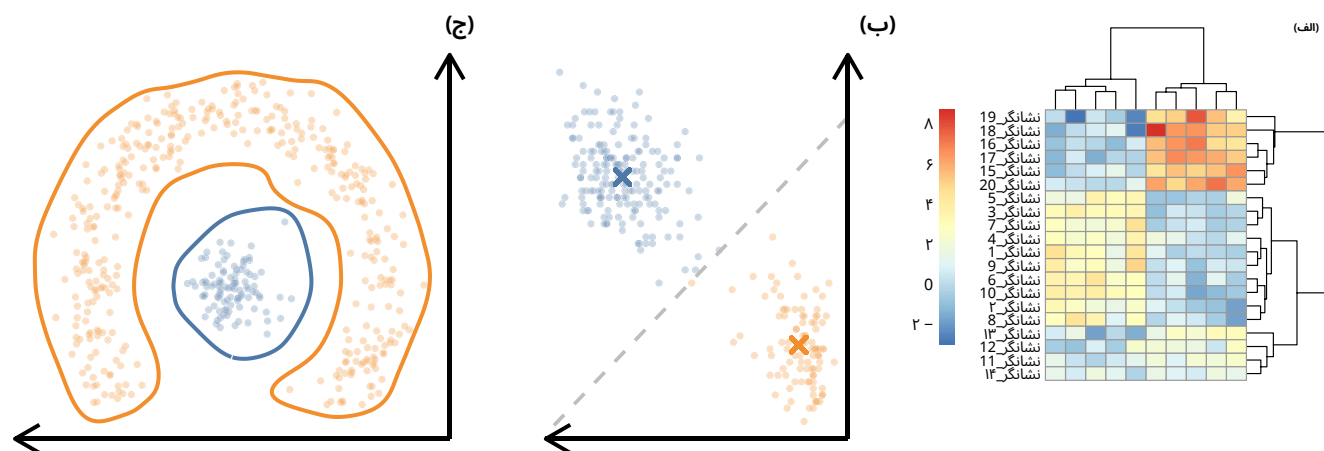
شناسایی خوشه های معنادار، که می تواند به ویژه با توجه به فضای ویژگی با ابعاد بالا چالش برانگیز باشد.

گروه دیگری از روش های خوشه بندی، خوشه بندی مبتنی بر چگالی است.^{۱۹} در روش های مبتنی بر چگالی، یک خوشه بخشی از فضای ویژگی را نشان می دهد که در آن نقاط داده متراکم هستند. نقاط داده متعلق به مناطقی از فضای ویژگی با چگالی کم، نویز در نظر گرفته می شوند. یکی از الگوریتم های خوشه بندی مبتنی بر چگالی شناخته شده، خوشه بندی مکانی مبتنی بر چگالی بر پایه های دارای نویز است.^{۲۰} خوشه بندی مبتنی بر چگالی نیازی به مقدار از پیش تعریف شده ای برای تعیین تعداد خوشه ها ندارد و نتیجه ای تکرارپذیر ارائه می دهد. علاوه بر این، قادر به شناسایی اشکال پیچیده خوشه ها، مانند آنچه در شکل نشان داده شده است، نیز می باشد. **شکل ۲ ج.**

در تحلیل خوشه بندی سلسله مراتبی، هدف ساخت سلسله مراتبی از خوشه ها است (ESL، فصل ۱۴).^۹ یک رویکرد ساده برای خوشه بندی سلسله مراتبی، اتصال همسایه است. ابتدا، تمام فواصل جفتی بین تمام نقاط داده در مجموعه داده ها محاسبه می شوند. سپس، در هر مرحله از یک فرآیند تکراری، دو نقطه داده با کمترین فاصله با هم گروه بندی می شوند. این منجر به یک ساختار خوشه بندی درخت مانند می شود، همانطور که در شکل نشان داده شده است. **شکل ۲ الف** در سمت چپ و بالای نقشه حرارتی که طول شاخه های درخت نشان دهنده فواصل نمونه ها است. برای رسیدن به مجموعه ای گسسته از خوشه ها مانند یک باید یک آستانه فاصله انتخاب شود که در آن درخت به صورت افقی بریده شود. باز هم، هیچ راه بهینه ای برای انتخاب چنین آستانه ای وجود ندارد و ممکن است راه حل های معقول زیادی وجود داشته باشد. خوشه بندی سلسله مراتبی می تواند به تنهایی یا در ترکیب با نقشه های حرارتی (مثلاً **شکل ۲ الف**) برای تجسم ویژگی های انتخاب شده یا همه آنها، به عنوان مثال، داده های بیان ژن.

کاهش ابعاد

تعداد ویژگی ها و بنابراین، ابعاد فضای ویژگی می تواند بسیار زیاد باشد و ده ها هزار معیار در هر نمونه داشته باشد. این امر نه تنها تجسم داده ها را چالش برانگیز می کند، بلکه تجزیه و تحلیل آنها را نیز چالش برانگیز می کند. به طور خاص، تجزیه و تحلیل مجموعه داده های با ابعاد بالا می تواند با پدیده ای به نام «نفرین ابعاد» مرتبط باشد.^{۲۱} که به پراکندگی داده ها و ویژگی های هندسی غیرشهودی در داده های با ابعاد بالا اشاره دارد



شکل ۲ مروری بر نتایج رویکردهای مختلف خوشه بندی. (الف) نتایج یک خوشه بندی سلسله مراتبی دوبعدی را نشان می دهد. دو دندروگرام، شباهت بین نمونه ها و همچنین بین نشانگرهای اندازه گیری شده را تجسم می کنند. چنین تجسمی اغلب در زیست شناسی برای بیان ژن یا سایر خوانش های فناوری omics استفاده می شود. (ب) نتیجه یک خوشه بندی کلاسیک را با استفاده از نشان می دهد. یعنی با مقدار انتخابی $k=2$. خوشه های حاصل معمولاً محدب هستند و هر نقطه به یک خوشه اختصاص داده می شود، یعنی خوشه ای که با نزدیکترین نقطه مرکزی (با X مشخص شده است) نشان داده شده است. (ج) نتیجه خوشه بندی مبتنی بر چگالی را نشان می دهد. لطفاً توجه داشته باشید که این رویکرد می تواند اشکال خوشه های غیرمحدب، مانند خوشه نارنجی، را شناسایی کند.

خروجی برای داده های جدید دیده نشده، که در آن مقادیر ورودی را مشاهده کرده ایم اما خروجی مرتبط با آنها را ندیده ایم.

دودسته اصلی از یادگیری نظارت شده وجود دارد: (۱) طبقه بندی که در آن مقادیر خروجی به صورت دسته بندی هستند، و (۲) رگرسیون که در آن مقادیر خروجی به صورت عددی هستند.

در بخش های بعدی، زمینه برازش مدل در یادگیری نظارت شده و مسئله رایج پیش برازش معرفی می شوند. سپس، توضیح می دهیم که چگونه عملکرد برای ابزارهای طبقه بندی و رگرسیون ارزیابی می شود (یعنی چگونه کیفیت نگاشت از ورودی ها به خروجی ها توسط الگوریتم ارزیابی شود). این جنبه ضروری است، زیرا مزیت اتخاذ روش های یادگیری ماشین اغلب حول محور دستیابی به عملکرد بالاتر با بده بستان تفسیرپذیری متمرکز است. درک معیارهای مختلف عملکرد، ارزیابی بهتر مزایای یک مدل پیشنهادی را امکان پذیر می کند، برخلاف این فرض که یک راه حل یادگیری ماشین همیشه می تواند از یک رویکرد سنتی بهتر عمل کند.

سپس به بررسی برخی از روش های طبقه بندی و رگرسیون موجود می پردازیم، از سطح پایین شروع می کنیم، جایی که تفسیر مدل ها هنوز سراسر است، و به سمت رویکردهای یادگیری ماشین محور ترپیش می رویم که در آن ها عملکرد، اغلب به قیمت از دست رفتن تفسیرپذیری، پیروز می شود. **شکل ۱** اطیف روش های موجود را از نظر عملکرد و قابلیت تفسیر خلاصه می کند. این بخش با مروری نه چندان جامع بر کاربردهای روش های یادگیری نظارت شده در زیست شناسی و به ویژه داروشناسی بالینی به پایان می رسد.

معیارهای عملکرد و مسئله ی بیش برازش

هدف یک الگوریتم یادگیری، یادگیری یک مفهوم یا تابع (= یک مدل) است که داده های آموزشی مشاهده شده را توصیف می کند و قادر است با اجتناب از کم برازش (underfitting) و بیش برازش (overfitting)، آن را روی داده های مستقل جدید تعمیم دهد.

عملکرد یک مدل با روش هایی ارزیابی می شود که امکان ارزیابی مدل را فراهم می کنند (یعنی تخمین میزان عملکرد کلی یک مدل مشخص و انتخاب مدل؛ و تخمین عملکرد مدل های مختلف برای انتخاب مناسب ترین مدل). برخی از این روش ها در بخش های بعدی برجسته شده اند.

برازش مدل، پارامترهای مدل بر اساس داده های مشاهده شده در مجموعه آموزشی تخمین زده می شوند. برای استخراج مقادیر بهینه پارامترها (مثلاً برای ضرایب و وزن ها)، یک معیار فاصله بین مدل و داده ها تعریف و به صورت عددی کمینه می شود. مستقل از معیار انتخاب شده، هدف از برازش مدل همیشه تخمین پارامترها با کمینه کردن فاصله است که به آن فاصله نیز گفته می شود. **تابع زیان** تابع هزینه، با دو الزام:

- مدل باید مقادیر پیش بینی شده ای را ارائه دهد که به مقادیر مشاهده شده در مجموعه آموزش نزدیک باشند، در غیر این صورت می گوئیم که **زیرپوش** ها و از درجه بالایی برخوردار است **تعصب**.
- مدل باید فراتر از مجموعه آموزشی تعمیم یابد. مدلی که **بیش برازش** روی مجموعه آموزش خوب پیش بینی می کند اما روی یک مجموعه آزمون مستقل ضعیف است، اغلب به این دلیل که برای داده ها خیلی پیچیده است. در این مورد، ما همچنین در **موردواریناس بالا**.

در ادامه، تماس خواهیم گرفت **تابع هدف** هر تابعی که برای تخمین پارامترهای مدل بهینه شده باشد.

فضاها. «نفرین ابعاد» چالش هایی را در اکثر رویکردهای تحلیل داده ها، از جمله یادگیری ماشین (ML) ایجاد می کند، اما محدود به آنها نیست.

برای کاهش چنین مشکلاتی، می توان از روش های کاهش ابعاد استفاده کرد. کاهش ابعاد می تواند با تبدیل هر نقطه داده با ابعاد بالا به دو یا چند بعد، ضمن حفظ اکثر تغییرات و فواصل نسبی، به تجسم داده ها کمک کند. علاوه بر این، حذف ویژگی های غیرمفید می تواند عملکرد مدل و زمان همگرایی را بهبود بخشد. اگرچه برخی از این روش ها، مانند تحلیل مؤلفه های اصلی، مدت ها قبل از ابداع اصطلاح یادگیری ماشین توسعه یافته اند، ۲۲ موارد دیگر، مانند جاسازی همسایه تصادفی توزیع شده با ۳۴ تا تقریب و تصویر منیفولد یکنواخت، ۱۳۴ اخیراً توسعه یافته اند و به چالش های پیچیده ای که در تجزیه و تحلیل داده ها ایجاد می شوند، می پردازند. همچنین یک رویکرد قدرتمند کاهش ابعاد مبتنی بر شبکه عصبی به نام خودرمزگذار وجود دارد. برای جزئیات بیشتر در مورد نحوه اعمال کاهش ابعاد در داده های زیست پزشکی، خواننده را به یک بررسی اخیر ارجاع می دهیم. ۲۵

نمونه هایی از کاربردهای یادگیری ماشینی بدون نظارت

خوشه بندی به طور گسترده در تجزیه و تحلیل داده های با ابعاد بالا، مانند آزمایش های رونویسی، متابولومیک و پروتئومیک، مورد استفاده قرار می گیرد. معمولاً از خوشه بندی سلسله مراتبی برای شناسایی عوامل اصلی مؤثر بر خوانش ها و همچنین برای شناسایی ماژول هایی با درجه بالایی از تنظیم همسو استفاده می شود. در توالی یابی تک سلولی، از خوشه بندی غیرسلسله ای برای درک انواع سلول های موجود در نمونه استفاده می شود. خوشه بندی همچنین برای شناسایی روابط بین بیماران، بافت ها، بیماری ها یا حتی علائم بیماری استفاده می شود. ۲۶-۲۹ خود ترکیبات دارویی نیز ممکن است بر اساس بیان ژن، حساسیت و خواص پروتئینی هدف دسته بندی شوند. ۳۰-۳۲ با هدف هدایت کشف دارو.

کاهش ابعاد به طور معمول در آزمایش های ترانسکریپتومیک و سایر آزمایش های omics-، معمولاً برای شناسایی داده های پرت و اثرات دسته ای بالقوه، استفاده می شود. در توالی یابی تک سلولی، تقریب منیفولد یکنواخت و تصویرسازی یا جاسازی همسایه تصادفی توزیع شده هم برای تجسم داده ها و هم برای خوشه بندی بعدی استفاده می شوند. ۳۴ کاهش ابعاد همچنین برای تجسم فضای شیمیایی با ابعاد بالا استفاده می شود. ۳۳ با عنوان یک مرحله پیش پردازش برای بهبود عملکرد یک مدل یادگیری ماشین. ۳۴

غذاهای اصلی

- خوشه بندی می تواند برای درک ساختار داده ها با گروه بندی مشاهدات مشابه با هم مورد استفاده قرار گیرد.
- ک خوشه بندی به معنای ابزاری ساده اما قدرتمند است، با این حال، تعداد خوشه ها باید از قبل مشخص شود.
- روش های مبتنی بر چگالی به تعداد از پیش تعیین شده ای از خوشه ها نیاز ندارند امکان شناسایی الگوهای پیچیده در داده ها را فراهم می کنند.
- خوشه بندی سلسله مراتبی، یک نمای کلی از رابطه در سطوح مختلف ارائه می دهد.
- کاهش ابعاد نه تنها برای تجسم داده ها، بلکه برای حذف ویژگی های غیرمفید نیز استفاده می شود.

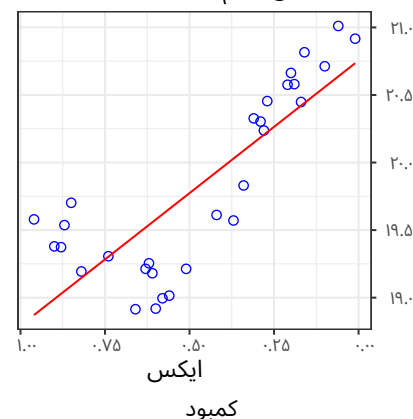
یادگیری تحت نظارت

دریک مسئله یادگیری نظارت شده، داده های آموزشی به همراه مشاهدات و مقادیر خروجی شناخته شده ی متناظر به کامپیوتر داده می شود. هدف، یادگیری قوانین کلی (که اغلب «مدل» نیز نامیده می شود) است که ورودی ها را به خروجی هانگاشت می کنند، به طوری که پیش بینی ... امکان پذیر باشد.

برازش چندجمله ای درجه ۱

خطای آموزش: ۰.۴

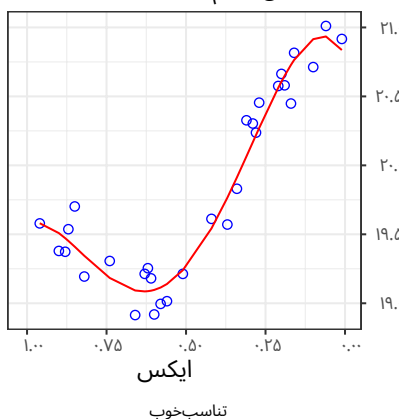
خطای تعمیم: ۰.۴۲



برازش چندجمله ای درجه ۴

خطای آموزش: ۰.۱۴

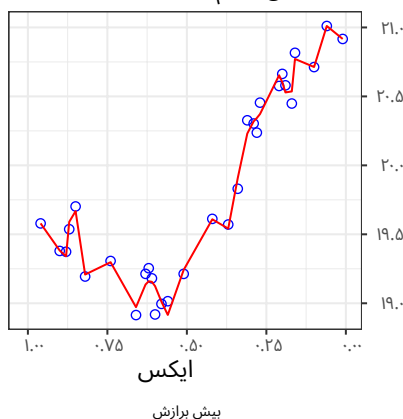
خطای تعمیم: ۰.۱۷



درجه برازش چندجمله ای ۲۰

خطای آموزش: ۰.۰۷

خطای تعمیم: ۲.۰۰



شکل ۳ نمایش مسئله ی کم برازش/بیش برازش در یک مورد رگرسیون ساده. نقاط داده با نقاط آبی و برازش مدل با خطوط قرمز نشان داده شده اند. کم برازش با یک مدل خطی (پنل سمت چپ)، برازش خوب با یک چندجمله ای درجه ۴ (پنل مرکزی) و بیش برازش با چندجمله ای درجه ۲۰ (پنل سمت راست) رخ می دهد. جذر میانگین مربعات خطابه عنوان تابع هدف برای ارزیابی خطای آموزش و خطای تعمیم انتخاب شده است که با استفاده از اعتبارسنجی متقابل fold-10 ارزیابی می شود.

در حالت رگرسیون، **شکل ۳** مسئله ی کم برازش و بیش برازش را در زمینه ی رگرسیون نشان می دهد. کم برازش می تواند زمانی رخ دهد که مدل خیلی ساده باشد یا وقتی ویژگی های استخراج شده از داده ها به اندازه ی کافی آموخته نشده باشند (**شکل ۳**، پنل سمت چپ). بیش برازش اغلب زمانی رخ می دهد که مدل بیش از حد پیچیده باشد یا ویژگی های زیادی در یک مجموعه کوچک از مثال های آموزشی وجود داشته باشد (**شکل ۳**، پنل سمت راست). این مشکل کم برازش/بیش برازش اغلب به عنوان بده بستان بایاس/واریانس نیز شناخته می شود، که از بیان خطای پیش بینی مورد انتظار، شامل هر دو عبارت بایاس و واریانس، ناشی می شود. بایاس نشانه ای از میانگین خطای مدل برای مجموعه های آموزشی مختلف است؛ این اختلاف بین میانگین مقادیر پیش بینی شده و میانگین واقعی است که ماسه ی در پیش بینی آن داریم. واریانس نشان دهنده حساسیت مدل به مجموعه آموزشی است؛ برای یک نقطه معین، با پراکندگی مقادیر پیش بینی شده در اطراف میانگین آنها مطابقت دارد. برای به حداقل رساندن خطای پیش بینی شده، بین به حداقل رساندن بایاس و واریانس، بده بستانی وجود دارد: افزایش پیچیدگی مدل، بایاس را کاهش می دهد اما واریانس را افزایش می دهد. برای ساخت مدل های ساده تر، تکنیک های مختلفی وجود دارد که تحت عنوان منظم سازی خلاصه شده اند. اصل این روش شامل اصلاح تابع هدف با اضافه کردن عبارات جریمه ای است که بر تخمین پارامتر تأثیر می گذارند. منظم سازی L1 و L2 رایج ترین آنها هستند (ESL)، بخش های 3.4.1 و 3.4.2.

دسته بندی های مختلف توابع زیان، هدف متفاوت

توابعی را می توان برای اندازه گیری فاصله بین داده های مشاهده شده و مقادیر پیش بینی شده توسط مدل انتخاب کرد. برخی از معیارهای فاصله ای که در عمل استفاده می شوند را می توان به ... مرتبط کرد. احتمال/درستنمایی نشان می دهد که مشاهده داده های ما طبق مدل انتخاب شده چقدر محتمل است. رایج ترین کاربرد درستنمایی، یافتن پارامترهایی است که مدل را به طور بهینه با داده ها برازش می دهند (یعنی تخمین های پارامتر درستنمایی حداکثری). معمولاً لگاریتم منفی درستنمایی به حداقل می رسد و به عنوان تابع هدف در نظر گرفته می شود زیرا خواص عددی مطلوبی دارد. به طور مشابه، در معیارهای یادگیری ماشین، مانند

میانگین مربعات خطا، تابع هدف لجستیک یا آنتروپی متقاطع، برای یافتن پارامترهای بهینه یا ارزیابی برازش مدل استفاده می شوند. در عمل، محاسبه تحلیلی حداکثر درستنمایی یا حداقل زیان ممکن است امکان پذیر نباشد و اغلب لازم است از یک الگوریتم بهینه سازی عددی برای حل مسئله و یافتن بهترین مقادیر پارامتر استفاده شود. *گرادیان نزولی* چنین الگوریتمی است که در آن ابتدا یک تابع هدف تعریف می کنیم که می خواهیم آن را به حداقل برسانیم و سپس به طور تکراری مقادیر پارامترها را در جهتی که شدیدترین کاهش (مشتق مرتبه اول) تابع هدف رخ می دهد، به روزرسانی کنید تا زمانی که همگرایی به حداقل فاصله حاصل شود. در سناریوی تابع هدف غیرمحدب، موفقیت یافتن یک حداقل سراسری، برخلاف رسیدن به برخی حداقل های محلی، به انتخاب مجموعه اولیه مقادیر پارامتر، نرخ یادگیری (یعنی اندازه گام هر تکرار) و معیار همگرایی بستگی دارد. خواننده می تواند به مرجع ... مراجعه کند. ۳۵ برای جزئیات بیشتر در مورد فرآیندهای بهینه سازی محدب و غیرمحدب. نزول گرادیان تصادفی یک ترفند اضافی است که می تواند با نمونه برداری تصادفی از یک مجموعه داده آموزشی و جمع کردن فواصل در این زیرمجموعه از نقاط داده آموزشی برای تقریب تابع هدف، بهینه سازی را سرعت بیشتری بخشد.

اصول کلی انتخاب و ارزیابی مدل. مشکل

بیش برازش نشان می دهد که عملکرد مدل روی مجموعه داده آموزشی، شاخص خوبی برای عملکرد آن روی مجموعه داده جدید نیست. در ادامه اصول ارزیابی عملکرد مدل در یک محیط یادگیری^۷ نظارت شده را شرح خواهیم داد.

اصل کلی انتخاب مدل به شرح زیر است: وقتی داده های کافی وجود داشته باشد، آنها را به سه زیرمجموعه تقسیم می کنیم - مجموعه های آموزشی، اعتبارسنجی و آزمون. مجموعه آموزشی برای ساخت مدل های مختلف استفاده می شود، در حالی که مجموعه اعتبارسنجی متعاقباً برای انتخاب الگوریتم و در صورت نیاز، انتخاب ابرپارامترها استفاده می شود. سپس، مدلی که بهترین عملکرد را در مجموعه اعتبارسنجی دارد انتخاب می شود. در نهایت، مجموعه آزمون امکان ارزیابی خطای تعمیم را فراهم می کند که به آن خطای تعمیم نیز گفته می شود. خطای آزمون که خطای پیش بینی روی یک مجموعه داده آزمایشی است

درمقابل، معیار اطلاعات بیزی:

$$(3) \quad \text{آی سی} = \ln(n) - 2 \cdot \text{لینر} (-),$$

تعدادنقاط داده را در نظر می گیرند.

این رویکردهای انتخاب مدل به ندرت در یادگیری ماشینی استفاده می شوند، که بخشی از آن به دلیل پیچیدگی مجموعه داده ها و نقض فرضیات توزیعی مرتبط است. در عوض، رویکردهایی مانند اعتبارسنجی متقابل (خوشه بندی) بیشتر مورد استفاده قرار می گیرند.

معیارهای عملکرد برای ارزیابی مدل برای رگرسیون در مدل ها، ما معمولاً از میانگین مربعات خطا یا انواع دیگر توابع هدف میانگین برای مقایسه عملکرد مدل در مجموعه آموزش و آزمون استفاده می کنیم. برای مسأله طبقه بندی دو کلاسه، معیارهای عملکرد رایج اغلب از «ماتریس سردرگمی» نشان داده شده در شکل ۵ به طور خلاصه در زیر توضیح داده شده است.

- دقت، مربوط به نسبت مقادیر مثبت درست پیش بینی شده به تعداد کل مقادیر مثبت پیش بینی شده.
- یادآوری، که نرخ مثبت واقعی (TPR) نیز نامیده می شود و مربوط به نسبت مقادیر مثبت پیش بینی شده ی صحیح به تعداد کل مقادیر مثبت در مجموعه داده ها است.
- نرخ مثبت کاذب (FPR)، مربوط به نسبت مقادیر منفی پیش بینی شده نادرست است.
- دقت، مربوط به تعداد مقادیر درست پیش بینی شده تقسیم بر تعداد کل مقادیر پیش بینی شده.
- مساحت زیر منحنی (AUC) ROC: منحنی های مشخصه عملیاتی گیرنده (ROC) وابستگی TPR (فراخوان) و FPR را نشان می دهند. چگالی در طبقه بندی دودویی، هر نقطه روی منحنی ROC با انتخاب آستانه های مختلف برای طبقه بندی ... مکان یابی می شود. من در کلاس مثبت یا منفی. گوشه بالا سمت چپ منحنی ROC

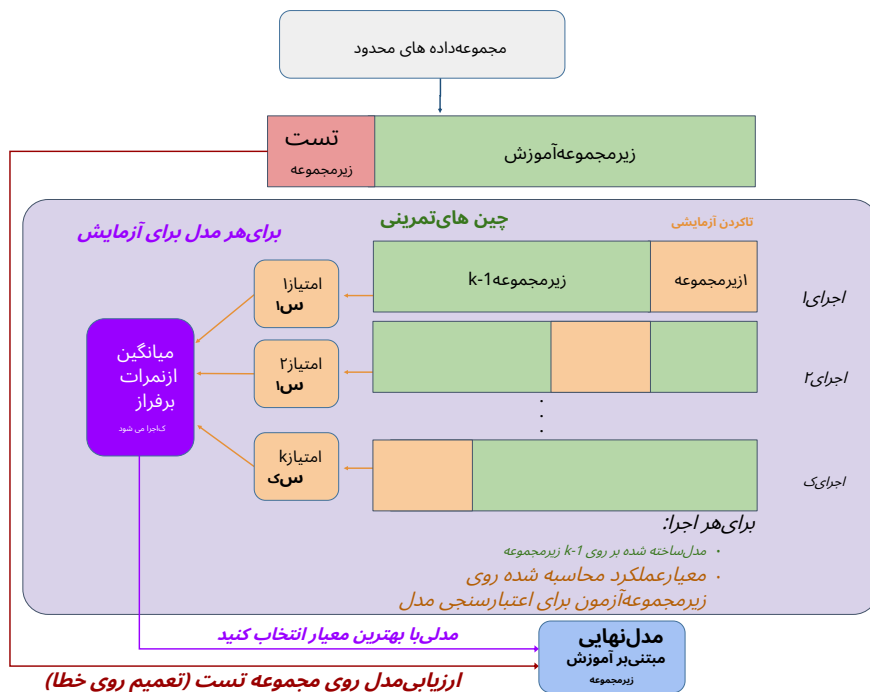
که در طول آموزش استفاده نشده است. لازم به ذکر است که وقتی مجموعه داده های اصلی دارای سوگیری باشند، خطای تعمیم می تواند بیشتر از حد انتظار باشد (به بخش «داده ها و ویژگی ها» مراجعه کنید). اعتبارسنجی مدل در برابر یک مجموعه داده آزمایشی کاملاً مستقل، روش استاندارد طلایی برای ارزیابی تعمیم پذیری مدل است.

وقتی مجموعه داده ها برای استخراج یک مجموعه اعتبارسنجی مناسب بسیار کوچک است، به عنوان مثال، می توان از تکنیک های اعتبارسنجی متقابل برای انتخاب ابرپارامترهای مدل استفاده کرد. پس از کنار گذاشتن زیرمجموعه ای از داده ها برای آزمایش، اعتبارسنجی fold شامل تقسیم مجموعه آموزشی به ... است. یک زیرمجموعه ها، یک-زیرمجموعه هایی که برای آموزش استفاده می شوند و زیرمجموعه آخری که برای ارزیابی عملکرد استفاده می شود. این فرآیند تکرار می شود بار، هر کدام یک زیرمجموعه یک بار برای اعتبارسنجی استفاده می شود و سپس نمرات عملکرد هر زیرمجموعه برای هر مجموعه از ابرپارامترها برای آزمایش، میانگین گیری می شوند. ک روش اعتبارسنجی متقابل برابر در خلاصه شده است. **شکل ۴** برای انتخاب بین الگوریتم های یادگیری مختلف ۳۰ می توان از اعتبارسنجی متقابل تو در تو استفاده کرد.

شاخص های پیچیدگی مدل در مقابل خوبی برازش. در فرارماکو متریک، انتخاب مدل معمولاً بر اساس معیارهای کمی است که خلاصه می کنند مدل چقدر خوب با داده ها برازش دارد، که اغلب با جریمه هایی برای بیش برازش همراه است. رایج ترین آنها معیار اطلاعات آکائی و معیار اطلاعات بیزی هستند. آنها تعداد پارامترهای مدل را جریمه می کنند و به خوبی برازش، که از طریق درستنمایی اندازه گیری می شود، پاداش می دهند. معیار اطلاعات آکائی یک به صورت زیر فرموله می شود:

$$(2) \quad \text{آی سی} = 2 - \ln(2) \cdot \text{لینر} (-),$$

با تعداد پارامترها و حداکثر احتمال .



شکل ۴ نمایش اصول کلی یادگیری نظارت شده در مورد یک مجموعه داده محدود. برای ارزیابی توانایی تعمیم یک الگوریتم یادگیری نظارت شده، داده ها به یک زیرمجموعه آموزشی که برای ساخت مدل استفاده می شود و یک زیرمجموعه آزمایشی که برای ارزیابی خطای تعمیم استفاده می شود، تقسیم می شوند.

برچسب های پیش بینی شده			
	0	1	
1	منفی کاذب (FN)	موقعیت واقعی (TP)	برچسب های واقعی (مشاهدات)
یادآوری=TPR (نرخ مثبت واقعی) $\frac{TP}{TP+FN}$ تنبی آر = $\frac{TP}{TP+FP}$			
0	منفی واقعی (TN)	موقعیت کاذب (FP)	
اختصاصی بودن = $\frac{TN}{TN+FP}$ نرخ مثبت کاذب: $\frac{FP}{TN+FP}$ برجسته سازی = $\frac{TP}{TP+FP}$			
دقت $\frac{TP+TN}{TP+FP+TN+FN}$	نرخ منفی کاذب $\frac{FN}{FN+TN}$	دقت $\frac{TP}{TP+FP}$	

شکل ۵: ماتریس درهم ریختگی برای مسائل دو کلاسه. ماتریس درهم ریختگی نشان می دهد که الگوریتم در پیش بینی برچسب ها در یک مسئله طبقه بندی دودویی که برچسب های مقادیر 0 (به نام «منفی») یا 1 (به نام «مثبت») را می گیرند، با ارزیابی برچسب های پیش بینی شده در مقابل برچسب های واقعی، چقدر موفق بوده است. هر نقطه داده در مجموعه آزمون متعلق به یکی از چهار دسته است و می توان معیارهای مختلفی را از این اعداد استخراج کرد.

یک روش طبقه بندی بسیار مؤثر در عمل، این روش در آموزش (یعنی ذخیره داده هادر پایگاه داده) و پیاده سازی های کارآمد برای محاسبه ی ... بسیار کارآمد است. ک شبکه های عصبی (NN) وجود دارند. خب، چالش های این رویکرد چیست؟ واضح ترین چالش این است که چون هیچ «مرحله یادگیری» وجود ندارد، ک رویکرد شبکه های عصبی، ویژگی هایی را که واقعاً برای پیش بینی کلاس یک مورد جدید مرتبط هستند، شناسایی نمی کند. بنابراین، حتی اگر در یک فضای ویژگی ۲۰ بعدی، که در آن فقط ۲ مورد ممکن است واقعاً برای طبقه بندی مرتبط باشند، فاصله با در نظر گرفتن هر ۲۰ بعد محاسبه خواهد شد. بنابراین، ک نزدیکترین نقاط داده ای که توسط پرس وجو برگردانده می شوند، به شدت تحت تأثیر ویژگی های نامربوط یا نویز قرار خواهند گرفت (همچنین به «کاهش ابعاد» در مورد نحوه حذف برخی از این ویژگی ها مراجعه کنید). در نتیجه، طبقه بندی حاصل به جای الگوی واقعی زیربنایی در داده ها، توسط نویز هدایت می شود. از این نظر، این رویکرد از همان چالشی رنج می برد که رویکردهای خوشه بندی (به بخش «خوشه بندی» مراجعه کنید) با آن مواجه هستند، که اغلب به عنوان «نفرین ابعاد» خلاصه می شوند. ۲۱.

بیز ساده

دومین رویکرد یادگیری بسیار شهودی که می خواهیم معرفی کنیم، بیز ساده است. این رویکرد مبتنی بر محاسبه آمار ساده از یک مجموعه داده آموزشی معین به عنوان مرحله یادگیری است که پس از یک کاربرد ساده (اما ساده) از فرمول بیزی برای احتمال شرطی به منظور دستیابی به یک طبقه بندی انجام می شود. به دلیل سادگی آن، اغلب برای دستیابی به یک عملکرد طبقه بندی پایه نیز استفاده می شود که سایر روش های پیچیده تر باید آن را بهبود بخشند. این رویکرد را می توان با یک مثال ساده به بهترین شکل توضیح داد. فرض کنید مجموعه داده های آموزشی شامل بیمارانی داریم که از سرماخوردگی بی ضرر یا عفونت آنفولانزا رنج می برند. ما برای هر بیمار دو ویژگی، یعنی تب (بالا، پایین یا بدون تب) و درد (قوی، پایین یا بدون تب) را اندازه گیری کرده ایم. برای هر بیمار، از طریق آزمایش آزمایشگاهی می دانیم که آیا بیمار به عفونت آنفولانزا مبتلا بوده است یا خیر. اکنون می خواهیم از این داده ها یاد بگیریم و آن را برای تشخیص یک بیمار جدید (در جایی که هیچ آزمایش آزمایشگاهی در دسترس نداریم) با استفاده از رویکرد بیز ساده به کار ببریم. به عنوان یک مرحله یادگیری، برای هر مقدار ویژگی، تعداد دفعات وقوع آن را در گروه بیمارانی مبتلا به آنفولانزا و سرماخوردگی محاسبه می کنیم (مثلاً برای به دست آوردن احتمال تب بالا در شرایط ابتلا بیمار به آنفولانزا و غیره). نتیجه این مرحله یادگیری را می توان در ... مشاهده کرد. جدول ۱.

حالت ایده آل این است که ۱۰۰٪ مقادیر مثبت به درستی طبقه بندی شوند ($TPR = 1$) و ۰٪ مقادیر مثبت به اشتباه در ۰ پیش بینی شوند ($FPR = 0$). از آنجایی که ایده آل این است که TPR را به حداکثر و FPR را به حداقل برسانیم، مساحت زیر منحنی (AUC) (ROC) بزرگتر بهتر است.

برخی از این معیارها را می توان برای مسائل چندکلاسه، که در آن ها بیش از دو برچسب مختلف در مجموعه داده ها وجود دارد، تعمیم داد. با این حال، معیارهای ذکر شده در بالا نسبت به پارامترهای مدل غیرپیوسته هستند، از این رو، بهینه سازی پارامتر ممکن است هنگام استفاده از آن ها به عنوان تابع هدف چالش برانگیز باشد. یک معیار جایگزین پیوسته و پرکاربرد که قبلاً در بخش «برازش مدل» ذکر شد، آنتروپی متقاطع است (ESL، فصل ۹)، که نه تنها محتمل ترین پیش بینی، بلکه امتیاز پیش بینی (اعتماد به پیش بینی) را نیز در نظر می گیرد.

ک-نزدیکترین همسایه ها

مابرسی اجمالی خود را بر روی روش های یادگیری موجود با روشی آغاز می کنیم که مرحله یادگیری را به طور کامل حذف می کند و بنابراین، منجر به یک مدل صریح که از داده های آموزشی آموخته می شود، نمی شود. همانطور که بعداً بحث خواهیم کرد، این نیز یکی از بزرگترین کاستی های آن است. این نوع یادگیری اغلب به عنوان «یادگیری مبتنی بر نمونه» نیز شناخته می شود و در مثال خاص ما، «ک-یادگیری نزدیکترین همسایه kNN » ۳۷. در این رویکردها، یادگیری صرفاً شامل ذخیره تمام نقاط داده موجود و برچسب گذاری شده (یعنی داده های آموزشی) در یک پایگاه داده است. هنگامی که یک مثال جدید، اما طبقه بندی نشده، مشاهده می شود، الگوریتم آن را در ... قرار می دهد. فضای ویژگی چندبعدی بر اساس مقادیر ویژگی آن. برای هر نقطه داده در پایگاه داده، اکنون فاصله (مثلاً فاصله اقلیدسی یا فاصله های پیچیده تر دیگر) را تا این نقطه داده جدید محاسبه می کنیم تا آن را شناسایی کنیم. ک-نزدیکترین همسایگان. در مرحله دوم، برچسب های شناخته شده این موارد را بررسی می کنیم که شبکه های عصبی در پایگاه داده ما. فرض کنید ما انتخاب کرده ایم اگر نه باشد، مشاهده می کنیم که هفت تا از نزدیکترین همسایه ها به عنوان کلاس X برچسب گذاری شده اند در حالی که دو تا از آنها به عنوان کلاس Y برچسب گذاری شده اند. در این حالت، نقطه داده جدید خود را به کلاس X اختصاص می دهیم زیرا اکثر همسایه های آن از این کلاس هستند. بسط این رویکرد ساده، وزن دهی به اهمیت همسایه ها در طبقه بندی بر اساس فاصله آنها تا نقطه داده جدید است. با وجود اینکه بسیار سراسر است و ساده است، اما ثابت می شود که ...

جدول 1 تصویرسازی بیز ساده: نمونه ای از نتایج مرحله یادگیری روی مجموعه داده آنفولانزا، که احتمال مقادیر ویژگی ها را با توجه به دسته بیمار نشان می دهد.

ویژگی ها			تب			درد		
کلاس ها			بالا			کم		
آنفولانزا (آنفولانزا)			آنفولانزا			آنفولانزا		
پ (تب = بالا آنفولانزا)	پ (تب = پایین آنفولانزا)	پ (تب = خیر آنفولانزا)	پ (درد = شدید آنفولانزا)	پ (درد = کم آنفولانزا)	پ (درد = بدون آنفولانزا)	پ (تب = بالا سرما)	پ (تب = پایین سرماخوردگی)	پ (تب = خیر سرماخوردگی)
0.1 =	0.95 =	0 =	0.75 =	0.20 =	0.05 =	0.1 =	0.4 =	0.5 =
سرد			سرد			سرد		
پ (تب = بالا سرما)	پ (تب = پایین سرماخوردگی)	پ (تب = خیر سرماخوردگی)	پ (درد = شدید سرما)	پ (درد = کم سرما)	پ (درد = بدون سرما)	پ (تب = بالا سرما)	پ (تب = پایین سرماخوردگی)	پ (تب = خیر سرماخوردگی)
0.9 =	0.4 =	0.5 =	0.3 =	0.3 =	0.4 =	0.1 =	0.4 =	0.5 =

جدول ۱ احتمال هر ویژگی را با توجه به دسته بیمار خلاصه می کند و نشان می دهد که در کل جمعیت بیمار، احتمال ابتلای یک بیمار به عفونت آنفولانزا ۰.۱ است، در حالی که احتمال ابتلا به سرماخوردگی معمولی ۰.۹ است.

وقتی این مقادیر را تولید کردیم و بنابراین «مرحله یادگیری» را با تجزیه و تحلیل مجموعه داده های خود تکمیل کردیم، بیز ساده لوحانه یک فرض ساده لوحانه مطرح می کند، و آن این است که همه این ویژگی ها از نظر شرطی مستقل از یکدیگر هستند. در واقعیت، این به ندرت صادق است و روش های یادگیری بیزی پیشرفته تری وجود دارند که این فرض را در نظر نمی گیرند. با این حال، این فرض امکان کاربرد ساده ای از قضیه بیزی را فراهم می کند. برای جزئیات (یعنی فرمول ها) در مورد نحوه استخراج این طبقه بندی کننده، خواننده را به مطالب بیشتر برای مطالعه (ESL، فصل 6) ارجاع می دهیم. به طور خلاصه، احتمال یک برچسب خاص (آنفولانزا یا سرماخوردگی) برای یک آیتم آزمایشی جدید را می توان به صورت حاصل ضرب احتمالات ویژگی شرطی واحد (تب و درد) که برای نقطه داده مشاهده می شوند، در احتمال کلاس (آنفولانزا یا سرماخوردگی) محاسبه کرد. کلاسی که بیشترین احتمال پسین را دارد، به عنوان کلاس پیش بینی شده برای آیتم آزمایشی انتخاب می شود. با فرض اینکه یک فرد آزمایشی با تشخیص ناشناخته آنفولانزا یا سرماخوردگی داریم و می دانیم که این فرد با تب بالا و تب بالا مراجعه می کند. اگر سطح درد بالا باشد، احتمال ابتلا به آنفولانزا را به صورت زیر محاسبه می کنیم:

$$\begin{aligned} & \text{پ (تب = بالا | آنفولانزا)} \times \text{پ (درد = شدید | آنفولانزا)} \\ & \text{پ (آنفولانزا)} = 0.007125 \times 0.10075 \times 0.95 = 0.000675 \end{aligned} \quad (4)$$

به همین ترتیب، احتمال سرماخوردگی را به صورت زیر محاسبه می کنیم:

$$\begin{aligned} & \text{پ (تب = بالا | سرماخوردگی)} \times \text{پ (درد = شدید | سرما)} \\ & \text{پ (سرمد)} = 0.9003 \times 0.1 = 0.09003 \end{aligned} \quad (5)$$

برای بیماری که با تب بالا و درد عضلانی یا سردرد شدید به پزشک مراجعه می کند، این امر منجر به احتمال خلفی (غیر نرمال) برای عفونت آنفولانزا ۷.۱۲۵٪ و احتمال ۲.۷٪ برای سرماخوردگی معمولی می شود. بنابراین، بیمار بیشتر از سرماخوردگی از آنفولانزا رنج می برد.

بنابراین، بیز ساده لوحانه، از بسیاری جهات، چگونگی یادگیری انسان ها از تجربه را فرموله می کند.

درخت های تصمیم، جنگل های تصادفی و تقویت گرادیان درخت های تصمیم، بلوک های سازنده ی ضروری برای بسیاری از الگوریتم های یادگیری ماشین هستند و حداقل ۵۰ سال است که مورد استفاده قرار می گیرند. ۳۸،۳۹ ایده پشت درخت های تصمیم گیری بسیار شهودی است و به بهترین شکل به صورت بصری نمایش داده می شود (مثلاً شکل ۱ بسته به مسئله، گره های برگ درخت تصمیم دارای کلاس، احتمال یا پیوسته هستند.

مقادیر در صورت رگرسیون. در روزهای اولیه یادگیری ماشین، از درخت های تصمیم برای حل مسائل دارویی مانند دوزبندی، سم شناسی و تشخیص استفاده شده است. ۴۰-۴۲ اگرچه استفاده از درخت های تصمیم گیری شهودی است، اما سوال این است که چگونه می توان چنین درخت هایی را از داده های موجود ساخت. چند رویکرد معروف که ارزش ذکر کردن دارند عبارتند از: CART و ID3. ۴۴

در حال حاضر، درخت های تصمیم تقریباً هرگز در یادگیری ماشین به شکل اصلی خود استفاده نمی شوند. یکی از دلایل این امر این است که درخت های تصمیم مستعد بیش برآزش هستند. با این وجود، درخت های تصمیم به سنگ بنای دو رویکرد پرکاربرد تبدیل شده اند: جنگل های تصمیم تصادفی و چارچوب های تقویت گرادیان.

هم جنگل های تصمیم گیری تصادفی و هم تقویت گرادیان مبتنی بر درخت از مجموعه ای (گروهی) از درخت های تصمیم گیری آموزش دیده برای پیش بینی متغیر نتیجه استفاده می کنند. تفاوت اساسی بین تقویت گرادیان مبتنی بر درخت و جنگل های تصمیم گیری تصادفی در نحوه ایجاد درختان است. در مورد جنگل های تصادفی، الگوریتم صدها یا هزاران درخت تصمیم گیری عمیق ("پیش بینی کننده های قوی") می سازد. هر یک از این درخت ها احتمالاً بیش برآزش شده اند، با این حال، با ترکیب خروجی های چندین درخت می توانیم مشکل آموزش بیش از حد را حل کنیم. برعکس، در یک الگوریتم افزایش گرادیان، مانند XGBoost یا CatBoost، هر یک از درخت های یک درخت تصمیم گیری کم عمق ("پیش بینی کننده ضعیف") هستند و الگوریتم به طور تکراری با اضافه کردن درختان بیشتر و بیشتر، خطای طبقه بندی را در طول زمان کاهش می دهد.

امروزه، روش های تقویت گرادیان، عملکرد بسیار خوبی را هم در نشریات و هم در مسابقات یادگیری ماشین نشان می دهند. حتی بدون تنظیم های پیرامون، آنها معمولاً عملکرد عالی را با هزینه محاسباتی نسبتاً کم ارائه می دهند. ۱۱۱ طرف دیگر، جنگل های تصادفی معمولاً کمتر مستعد بیش برآزش هستند. ۴۵ و نیاز به تنظیم پارامترهای کمتری دارند. ۴۶ این امر، جنگل های تصمیم گیری تصادفی را برای مجموعه داده های کوچک تر یا به عنوان یک روش پایه برای الگوبرداری جذاب می کند.

روش های گروه درختی می توانند برای وظایف طبقه بندی و همچنین برای رگرسیون استفاده شوند. در هر دو مورد، خروجی های درخت میانگین گیری می شوند که می تواند یک تابع خروجی هموار ایجاد کند.

روش های هسته: ماشین های بردار پشتیبان و رگرسیون روش های کرنل و به طور خاص تر، ماشین بردار پشتیبان (SVM) برای طبقه بندی و رگرسیون بردار پشتیبان (SVR) برای خروجی پیوسته، به دلیل توانایی شان در مقاوم بودن در برابر نویز و کار با مجموعه داده های با ابعاد بالا که در ژنتیک، ترانسکریپتومیکس و پروتئومیکس یافت می شوند، کاربردهایی در زیست شناسی محاسباتی پیدا کرده اند. ۴۷ به طور مشخص در یک مثال جدیدتر، از SVR برای مشخص کردن ترکیبات سلولی از داده های حجیم رونویسی استفاده شد. ۴۸

فضای ویژگی. یک تابع هسته که در فضای ورودی اعمال می شود، معادل یک ضرب نقطه ای در فضای ویژگی است که در آن معیارهای شباهت محاسبه می شوند. این امر بدون نیاز به نگاشت صریح داده های ورودی از فضای ورودی به یک فضای ویژگی توسط یک تابع نگاشت، محقق می شود.

بادر دست داشتن تمام این مفاهیم، اکنون می توانیم مدلی با ضخامت مشخص، که به عنوان لوله ای که توسط تابع زیان insensitive- β معرفی می شود، برازش دهیم، در حالی که عبارت منظم سازی، مسطح بودن این ابرصفحه را در فضای ویژگی تعریف شده توسط تابع هسته کنترل می کند. **شکل ۶** این مفاهیم اساسی SVM را نشان می دهد.

ترفندو انتخاب هسته: می تواند توابع هدف غیرخطی را ثبت کند که ورودی های چند متغیره را به خروجی نگاشت می کنند. به طور دقیق تر، ترفند هسته به این معنی است که یک هسته SVR

$$(6) \quad K((\mathbf{x}_1, \Phi), (\mathbf{x}_2, \Phi)) = \exp(-\gamma \| \mathbf{x}_1 - \mathbf{x}_2 \|^2)$$

اعمال شده بر مجموعه ای از ورودی ها در فضای ورودی معادل محاسبه ضرب نقطه ای به عنوان معیار شباهت در برخی از فضای ویژگی ها است. این امر بدون نیاز به انجام صریح پیش نگاشت ورودی ها حاصل می شود. **ایکس من**، با یک تابع نگاشت، یک تابع هسته محاسبه شده در فضای ورودی، اگر و تنها اگر یک تابع متقارن مثبت معین باشد، معادل یک ضرب نقطه ای در یک فضای ویژگی است. Φ است. انتخاب هسته تابع پایه شعاعی،

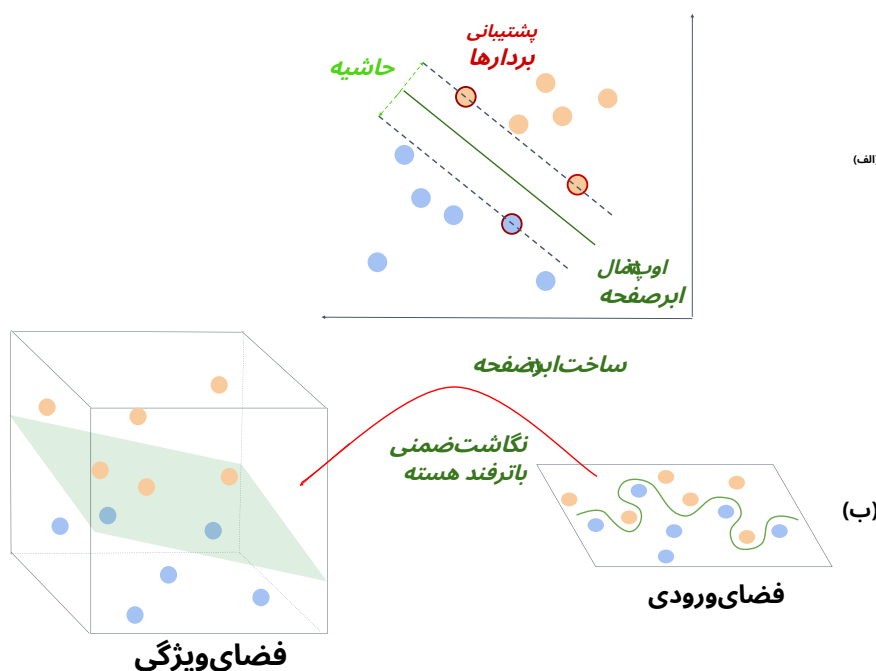
$$(7) \quad K((\mathbf{x}_1, \Phi), (\mathbf{x}_2, \Phi)) = \exp(-\gamma \| \mathbf{x}_1 - \mathbf{x}_2 \|^2)$$

اغلب به این دلیل ساخته می شود که می توان آن را به یک فضای ویژگی با ابعاد نامتناهی گسترش داد. اگرچه تابع پایه شعاعی طیف وسیعی را پوشش می دهد

این بخش ابتدا مروری مختصر بر مفاهیم کلیدی ارائه می دهد که مفاهیم تبدیلات هسته، یک تابع هدف با ناحیه بدون اتلاف و یک عبارت منظم سازی را برجسته می کند. تأکید بر ارائه استدلال پشت این موضوع خواهد بود که چرا این روش در برخورد با ورودی های چندگانه که اثرات آنها بر خروجی ناشناخته است و می توان فرض کرد که به توابع غیرخطی تبدیل می شوند، روشی تطبیق پذیرتر است.

پیشینه. مشابه تمام روش های رگرسیون، هدف SVR فرض تابعی روی ورودی (ها) است که می تواند به تخمین خروجی مشاهده شده کمک کند. به همین ترتیب، برای SVM، هدف یافتن مرز تصمیم گیری بهینه است که کلاس ها را از هم جدا می کند. همانطور که از نامش پیداست، مفهوم اصلی پشت SVR/رگرسیون، توانایی انتخاب عینی زیرمجموعه ای از داده های آموزشی به نام ... است. **بردارهای پشتیبانی** این بردارهای پشتیبان، مدل را تعریف می کنند که معمولاً یک ابرصفحه در فضای ویژگی است. برای دستیابی به این هدف، باید چندین مفهوم معرفی شوند.

- یک از دست دادن غیر حساس تابع اجازه می دهد تا باقیمانده های کمتر از، بدون اتلاف در نظر گرفته شوند و بنابراین، بخشی از بردارهای پشتیبان در نظر گرفته شده برای تخمین تابع خروجی-ورودی نباشند.
- الف اصطلاح منظم سازی با هدف جستجوی مدلی برای توصیف رابطه بین متغیرهای ورودی و خروجی به تابع هدف اضافه می شود، به طوری که ابرصفحه تا حد امکان مسطح نگه داشته شود.
- متغیرهای اسلک می توان برای در نظر گرفتن خطاهای آموزش، که حاشیه نرم نامیده می شود، زمانی که خروجی خارج از ناحیه غیر حساس θ یافت می شود، معرفی کرد. با معرفی متغیرهای slack، تحمل برای عبارت باقیمانده که بزرگتر از θ باشد، ایجاد می شود.
- الف تابع هسته به ما اجازه می دهد در فضایی با ابعاد بالاتر کار کنیم.



شکل ۶: تصویرسازی اصول ماشین بردار پشتیبان (SVM). (الف) تصویرسازی از یک حالت ساده که در آن یک ابرصفحه، دو گروه را مستقیماً در فضای ورودی ها از هم جدا می کند. (ب) نمایش انجام طبقه بندی غیرخطی با نگاشت ضمنی ورودی ها به فضاهای ویژگی با ابعاد بالا که در آن نقاط داده می توانند توسط یک ابرصفحه از هم جدا شوند.

شبکه عصبی بازگشتی. شبکه های عصبی بازگشتی، دسته ای از شبکه های عصبی هستند که به مجموعه داده های سری زمانی اختصاص دارند، زیرا رابطه تریبی ذاتی مشاهده شده در داده های یک نقطه زمانی با نقطه زمانی دیگر را در نظر می گیرند. این شبکه ها در چیزی که در این زمینه به عنوان ... شناخته می شود، موفقیت کسب کرده اند. داده های تریبی، جایی که ترتیب یا توالی زمانی سیگنال نقشی ایفا می کند، یعنی در پردازش زبان طبیعی و پیش بینی سری های زمانی. این روش که ارتباط نزدیکی با حوزه تحقیقاتی ما دارد، در پیش بینی نتایج حاصل از پرونده های الکترونیکی سلامت کاربرد پیدا کرده است، جایی که غنای آن ذاتاً از ساختار همبستگی توالی داده ها ناشی می شود تا اقدامات سریع و حتی پیش بینی شده ای را که باید توسط کادر پزشکی انجام شود، توصیه کند. ۵۷. طرح مجدد سوال برای حل یک معماری مدل سازی در حوزه فارماکومتری، تنها زمانی شروع به ظهور کرد که این مقاله در حال نگارش بود. تانگو و همکاران یکی از تلاش های نادر در مورد نحوه استفاده از ML (اینجا: RNN ها) برای توصیف PK رمیفانتیل را ارائه داده و نتایج را با روش استاندارد طلایی فارماکومتری NONMEM مقایسه کرده اند. ۵۸. اگرچه از مدل های غیراستاندارد PK برای مقایسه استفاده شد و تعمیم پذیری نتایج می تواند مورد چالش قرار گیرد، تانگو و همکاران... سهم ارزشمندی در نشان دادن کاربردهای RNN ها در فارماکومتری داشته باشد.

شکل پایه یک RNN در شکل نشان داده شده است. شکل ۷ ب که در آن هر حالت فعلی (در زمان t) با ترکیبی از حالت قبلی سیستم و ورودی فعلی تعریف می شود، که مشابه مفهوم سیستم های دینامیکی کلاسیک است. وزن های هر لایه را می توان تعیین کرد که چقدر باید به عقب نگاه کرد، مشابه یک ثابت زمانی. برخلاف شبکه عصبی پیش خور، یک وزن یکسان در بلوک واحد نوری منفرد در تمام مراحل زمانی گسسته قبلی به اشتراک گذاشته می شود.

در هسته RNN، از یک توالی ورودی تشکیل شده است که توسط ... تعریف می شود. ایکس (تی)، یک دنباله خروجی مطابق تعریف زیر (تی)، یک توالی حالت پنهان یاسیستی که به صورت زیر تعریف می شود (ت) و همچنین زیرمدول های زنجیری از واحدهای تکراری.

مراحل مورد نیاز برای آموزش یک مدل RNN به شرح زیر است:

۱. معماری شبکه را تعریف کنید و مدل را با وزن ها و بایاس های تصادفی مقداردهی اولیه کنید.
۲. برای محاسبه خروجی تخمینی، یک انتشار رو به جلو انجام دهید.
۳. خطا را در لایه خروجی محاسبه کنید.

از اثرات احتمالی، منجر به تفسیر دشوارتر مدل نهایی می شود. در عمل، انتخاب تابع هزینه بر اساس کارایی محاسباتی است. سایر هسته های محبوب شامل هسته های خطی و چندجمله ای هستند. ۴۷.

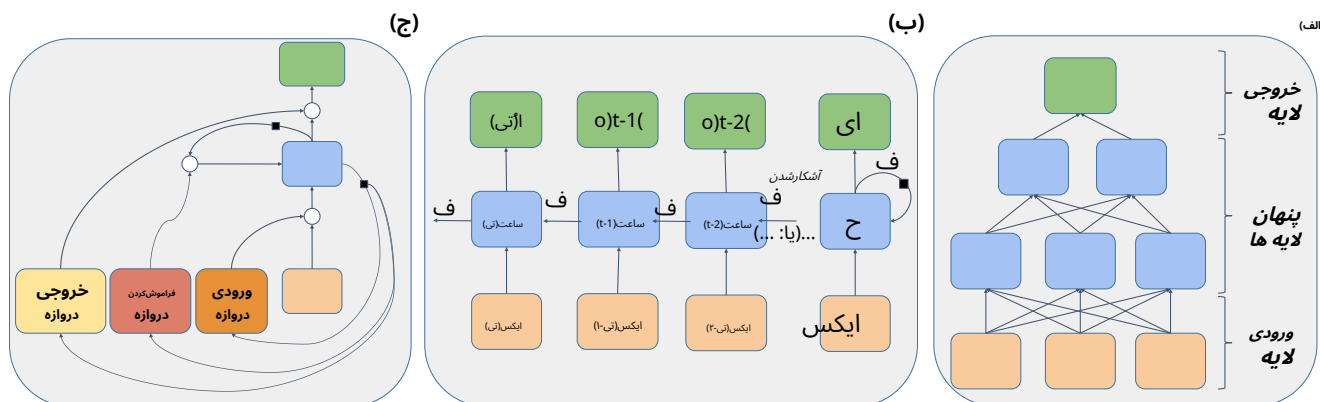
شبکه های عصبی

پیشینه. شبکه های عصبی مجموعه ای از نورون ها و لایه ها را تشکیل می دهند که منشأ آنها از تحلیل مدار گرفته شده است. می توان وزن های مختلفی را به هر لایه متصل کننده نورون ها اعمال کرد. در هر نورون، یک تابع فعال سازی به یک سیگنال ورودی وزن دار اعمال می شود تا یک سیگنال خروجی تولید شود. اغلب از یک تابع سیگموئیدی استفاده می شود که شامل یک فیلتر پایین گذر مرتبه اول از یک تابع پله واحد است. چنین تابع سیگموئیدی مزایای خروجی کران دار و مشتق پذیری پیوسته را دارد که در مرحله انتشار رو به عقب برای تنظیم وزن ها (پارامترهای مدل) مورد نیاز است، به مراحل تعریف شده در بخش «شبکه عصبی بازگشتی» مراجعه کنید.

نورون ها همانطور که در شکل نشان داده شده است، به لایه ورودی، لایه (های) پنهان و لایه خروجی تقسیم می شوند. شکل ۷ الف لایه های پنهان، لایه انتزاع مورد نیاز برای رفتن از لایه ورودی به لایه خروجی را انجام می دهند. تعداد لایه های پنهان مشخص می کند که آیا سیستم یک سیستم یادگیری سطحی (با یک یا چند لایه پنهان) یا یادگیری عمیق (با لایه های پنهان زیاد) است. بین تعداد لایه های پنهان و زمان مورد نیاز برای آموزش مدل، یک بده بستان ذاتی وجود دارد. به همین دلیل، اگرچه مفهوم اصلی تعبیه شده در شبکه عصبی، مفهوم جدیدی نیست، اما به دلیل پیشرفت های اخیر در قدرت محاسباتی، کاربردهای جدیدی پیدا کرده است.

اساسی ترین نوع آن به عنوان شناخته می شود پیشخور شبکه عصبی، زیرا اطلاعات فقط از لایه ورودی به لایه (های) پنهان و در نهایت به لایه خروجی منتشر می شود. وضعیت فعلی سیستم توسط هیچ وضعیت گذشته ای تعریف نمی شود؛ از این رو، یک سیستم بدون حافظه را نشان می دهد.

در ادامه، نمونه هایی از شبکه های عصبی شرح داده شده است: شبکه های عصبی بازگشتی، شبکه های حافظه کوتاه مدت بلندمدت و شبکه های بازگشتی دروازه دار. شبکه های عصبی قابل توجه دیگری که خارج از محدوده این مقاله هستند اما مطالعه بیشتر در مورد آنها را توصیه می کنیم، عبارتند از شبکه های عصبی کانولوشنی، ۵۳ شبکه های رمزگذار-رمزگشا، ۵۴ و مدل های تولیدی، ۵۶.



شکل ۷ شبکه های عصبی (الف) مبانی شبکه های عصبی پیش خور، (ب) آشکار شدن شبکه های عصبی بازگشتی، (ج) بسط شبکه های عصبی بازگشتی با واحدهای دروازه ای. مربع سیاه نشان دهنده تأخیر یک گام زمانی گسسته است.

۴. با استفاده از یک رویکرد بهینه سازی، یک الگوریتم انتشار معکوس برای به روزرسانی وزن ها انجام دهید.
۵. مراحل ۲ تا ۴ را برای تعداد دوره ها (یا تکرارها) تکرار کنید تا مقدار تابع زیان به حداقل برسد.

افزونه های این وانیل شبکه های عصبی بازگشتی (RNN) برای رسیدگی به مشکلات مربوط به گرادیان ناپایدار (مثلاً مشکل گرادیان ناپدیدشونده و معادل جدی تر آن، ناپایداری ناشی از گرادیان انفجاری) توسعه داده شدند. این مشکلات در اصل به دلیل ضرب هایی (تحت تأثیر خطاهای عددی) هستند که در انتشار معکوس در رابطه با تخمین خطا نسبت به پارامترها در امتداد هر لایه از شبکه عصبی ایجاد می شوند. به عبارت دیگر، گرادیان ناپدید شونده باعث می شود اطلاعاتی که باید از یک نقطه زمانی دورتر از زمان فعلی ثبت شوند، ضعیف شوند و بنابراین، مدل را برای ثبت حافظه ذخیره شده ارزشمند با تأخیر زمانی طولانی تر، ضعیف می کنند. در یک رویکرد کمتر رایج که حداقل یک مشتق جزئی الزام پایداری را نقض کند، و به ماتریس حالت داشتن حداقل یک مقدار ویژه $1 <$ تبدیل شود، این امر منجر به یک مشکل گرادیان انفجاری می شود، یک مشکل شناخته شده در سیستم پویای سنتی برای زمان گسسته. وظیفه مورد استفاده برای رسیدگی به این مشکل اساسی، در دو افزونه شناخته شده RNN (حافظه کوتاه مدت طولانی (LSTM) و شبکه بازگشتی دروازه دار (GRU)) بیشتر توضیح داده خواهد شد.

انواع توسعه های مختلفی در تحقیقات RNN وجود داشته است، هر روش جدید برای پرداختن به یک مسئله متفاوت به کار می رود که در نهایت منجر به توسعه مدل های قوی تر می شود. به عنوان مثال، برای دور زدن مشکل گرادیان ناپایدار، برش گرادیان با مجبور کردن گرادیان به یک آستانه در مرجع پیشنهاد شده است. ۹۰،۶۰، اما تاکنون پذیرفته شده ترین روش، گنجاندن واحدهای دروازه ای.

حافظه کوتاه مدت بلند و شبکه بازگشتی دروازه دار LSTM بخشی از خانواده بزرگتری از RNN های دروازه دار است که با معرفی واحدهای دروازه ای، اطلاعات را حفظ و فراموش می کنند. به طور خاص تر، همانطور که در نشان داده شده است، می توان سه واحد دروازه ای را در سیستم گنجانده. شکل ۷ ج اول، کپی کردن یا پاک کردن مستقیم کل حالت می تواند توسط گیت فراموشی کنترل شود. رویکرد مشابهی نیز توسط گیت ورودی برای تصمیم گیری در مورد اینکه آیا سیگنال ورودی فعلی را به عنوان بخشی از به روزرسانی حالت در نظر بگیرد یا خیر، انجام می شود. مقدار اطلاعاتی که باید از سیگنال حالت قبلی و از سیگنال ورودی اختلال حفظ شود، در هر گام زمانی آموخته می شود. سیستم باید با حفظ اطلاعات، وابستگی های زمانی بلندمدت را یاد بگیرد، اما باید گاهی اوقات نیز یاد بگیرد که اطلاعات را از وضعیت فعلی خود پاک کند. ۶۲. در نتیجه، حل مسائل گرادیان ناپدید شونده و انفجاری. در نهایت، می توان یک گیت خروجی، هرچند کمتر رایج، را به عنوان یک مکانیسم گیتینگ برای تصمیم گیری در مورد اینکه کدام سیگنال خروجی به سیستم بازگردانده شود، معرفی کرد.

یک تفسیر ساده تر و در نتیجه، پیاده سازی سریع تر آموزش را می توان در GRU یافت. GRU ها همان مشکل گرادیان های ناپایدار را برطرف می کنند و یک افزونه جدید به این خانواده از افزونه های RNN هستند. تفاوت اصلی بین LSTM و GRU این است که دومی گیت خروجی را حذف کرده و از گیت های ریست و به روزرسانی ساده تری استفاده می کند. ۶۳ این حال، در تئوری، LSTM باید عملکرد بهتری داشته باشد زیرا می تواند اطلاعات را از فاصله زمانی/تاخیر طولانی تر، وزن دهی یا وزن دهی کند.

نمونه هایی از کاربردهای یادگیری ماشینی تحت نظارت در داروشناسی بالینی

مدل های موجود در فارماکولوژی بالینی معمولاً با ترجمه اصول فیزیولوژیکی و دارویی به سیستم های معادلات دیفرانسیل و استفاده از الگوریتم های حداکثرسازی انتظار برای تخمین پارامترهای مدل ایجاد شده اند. این رویکرد با انگیزه مکانیکی در بسیاری از کاربردها مفید بوده و یک جزء تثبیت شده از برنامه های توسعه دارو است. به طور بالقوه به دلیل موفقیت این رویکردهای تثبیت شده، تاکنون تنها چند نمونه از اعمال روش های یادگیری ماشینی در مسائل فارماکولوژی بالینی وجود دارد. ریو همکاران یک شبکه عصبی عمیق رابر روی یک پایگاه داده بزرگ و منظم که شامل ۱۹۲۲۸۴ تداخل دارو با دارو بود، آموزش داد تا تداخلات دارو با دارو و دارو و غذا را برای نسخه ها، توصیه های غذایی و مولکول های جدید پیش بینی کند. ۶۴ ترکیب مجموعه داده های حاصل از مطالعات متعدد برای ایجاد پایگاه های داده بزرگ، پتانسیل استفاده از یادگیری ماشینی را برای پرداختن به سوالات گسترده داروشناسی بالینی افزایش می دهد.

همچنین یادگیری ماشینی برای ایجاد پل بین کشف دارو و توسعه بالینی استفاده شده است. به عنوان مثال، هامان و همکاران توانستند با استفاده از روش درخت تصمیم گیری، میزان بروز عوارض جانبی را از ساختار شیمیایی یک مولکول پیش بینی کنند. ۶۵ به طور مشابه، لنکستر و سوبی SVM ها را برای پیش بینی خطر de Pointes Torsades از ... پیاده سازی کردند. در شرایط آزمایشگاهی (in vitro) داده ها ۶۶

در حوزه ایمنی شخصی سازی شده، یادگیری ماشین توسط دانهاور مورد استفاده قرار گرفته است. و همکاران شخصی سازی ایمنی در زمینه هیپریلیبروبینمی در نوزادان. ۶۷ نویسنندگان از لاسو و جنگل های تصادفی برای پیش بینی از مجموعه داده های بالینی استفاده کردند. علاوه بر این، یادگیری تقویتی توسط گاودا مورد استفاده قرار گرفت. و همکاران برای شخصی سازی مدیریت دارویی کم خونی. ۶۸ رویکرد مشابهی برای توسعه یک سیستم «حلقه بسته» برای کنترل گلوکز با ترکیب یک مدل ریاضی، یک حسگر گلوکز و یک مدل یادگیری تقویتی استفاده شد. ۶۹ چاودا و همکاران و هنینگو همکاران امکان سنجی بازخورد بیزی برای تنظیم دوز آنتی بیوتیک ها را بررسی کرد. ۷۰، ۷۱ حوزه مراقبت های بهداشتی شخصی سازی شده می تواند از استفاده از مدل های یادگیری ماشین که تنظیمات دوز را در زمان واقعی توصیه می کنند، بسیار بهره مند شود. در یک مطالعه اخیر، یک الگوریتم کنترل از نوع یادگیری ماشین با مدل های ساختاری موجود PK/PD که برای متخصصان داروشناسی آشنا هستند، ادغام شد و مشخص شد که سیستم کنترل حلقه بسته حاصل، عملکرد بهتری نسبت به پمپ با کمک حسگر دارد. ۶۹

غذاهای اصلی

- روش های یادگیری نظارت شده، مدل ها را بر اساس جفت های خروجی-ورودی برچسب گذاری شده ی مجموعه داده های آموزشی استنتاج می کنند.
- معیارهای عملکرد برای ارزیابی مدل های طبقه بندی و رگرسیون استفاده می شوند تا از بیش برآزش مجموعه داده های آموزشی جلوگیری شود.
- روش های یادگیری نظارت شده ی زیادی با بده بستان های مختلف بین تفسیرپذیری و عملکرد وجود دارند.
- شکل خاصی از شبکه عصبی است که یک سیستم پویا را در زمان گسسته نمایش می دهد RNN •
- نمونه هایی از کاربردهای این روش های یادگیری تحت نظارت در زیست شناسی محاسباتی و به ویژه داروشناسی بالینی در حال پدیدار شدن هستند.

بحث

در این آموزش، ما برخی از روش های اساسی یادگیری ماشینی را معرفی کرده ایم که احتمالاً برای فارماکولوژی بالینی جالب خواهند بود.

توسعه به دلیل این موفقیت، علیرغم ورود یادگیری ماشینی، انتظار نمی رود که اهمیت و فعالیت رویکردهای کلاسیک فارماکومتری کاهش یابد. در مقابل، می توان آنها را با دانش و بینشی که توسط روش ها و مدل های یادگیری ماشینی استخراج شده است، تقویت و بهبود بخشید.

یک چالش مداوم برای اعضای جامعه داروشناسی بالینی که مایل به استفاده از روش های یادگیری ماشین هستند، شیوع ذاتی داده های طولی است. تاکنون، روش های یادگیری ماشین زیادی وجود دارند که برای پیش بینی به ویژگی های پایه متکی هستند، اما نمونه های نسبتاً کمی وجود دارد که در آن ها از داده های طولی استفاده شده باشد.

در مجموع، ما انتظار داریم که هرگز یک رویکرد جهانی و یکسان برای همه وجود نداشته باشد که مدل سازان از حوزه های مختلف به آن روی آورند. ما خاطرنشان می کنیم که حوزه های زیادی از هم افزایی بالقوه وجود دارد که در آن ها حوزه های مدل سازی در حوزه توسعه دارو با هم همپوشانی دارند. جامعه داروشناسی بالینی به پایه گذاری تحلیل های خود بر اساس اصول دارویی ادامه خواهد داد و به تدریج عناصر جدید یادگیری ماشین را به گردش کار خود اضافه خواهد کرد و مدل های خود را بیشتر تقویت خواهد کرد. علاوه بر این، جامعه داروشناسی بالینی قادر خواهد بود طیف وسیعی از سؤالاتی را که می تواند با استفاده از رویکردهای یادگیری ماشین به آن ها پاسخ دهد، افزایش دهد.

تأمین مالی

تأمین مالی شده است F. Hoffmann-La Roche Ltd. این مطالعه توسط Roche دریافت کننده یورسپه ی پسادکتری IID

تضاد منافع

همه نویسندگان هیچ گونه تضاد منافی برای این اثر اعلام نکردند.

© 2020 نویسندگان. فارماکولوژی بالینی و درمان شناسی منتشر شده توسط انتشارات وایلی، به نمایندگی از انجمن آمریکایی فارماکولوژی و درمان شناسی بالینی.

این یک مقاله با دسترسی آزاد تحت شرایط مجوز NonCommercial-NoDerivs Creative Commons Attribution است که استفاده و توزیع در هر رسانه ای را مجاز می داند. مشروط بر اینکه به اثر اصلی به درستی استناد شود، استفاده غیرتجاری باشد و هیچ گونه تغییر یا اقتباسی انجام نشود.

۱. کاماچو، دی ام، کالینز، کی ام، پاورز، آر کی، کاستلو، جی سی و کالینز، جی جی. یادگیری ماشینی نسل بعدی برای شبکه های بیولوژیکی. *سلول* ۱۵۸۱، ۱۵۹۲-۲۰۱۸.

۲. شن، دی، وو، جی. و سوک، اچ. آ. یادگیری عمیق در تحلیل تصاویر پزشکی. *سال. کشیش. زیست پزشکی. مهندسی* ۲۲۱، ۲۴۸-۲۰۱۷.

۳. راجکومار، آ.، دین، جی. و کوهان، آی. یادگیری ماشین در پزشکی.

مجله پزشکی انگلستان (N. Engl. J. Med) ۳۷۳، ۱۳۵۸-۲۰۱۹.

۴. کلاین، اس. سی. نمایش رویدادها در شبکه های عصبی و اتوماتای متناهی (نیروی هوایی پروژه رند، سانتا مونیکا، کالیفرنیا، ۱۹۵۱) <ADA596138/https://apps.dtic.mil/docs/citations>.

۵. بریمن، ل. مدل سازی آماری: دو فرهنگ (همراه با توضیحات و پاسخ نویسنده). *آمار علم* ۱۹۹، ۲۳۱-۲۰۰۱.

۶. ریبریو، ام تی، سینگ، اس. و گسترین، سی. «چرا باید به شما اعتماد کنم؟»: توضیح پیش بینی های هر طبقه بندی کننده مجموعه مقالات بیست و دومین کنفرانس بین المللی ACM SIGKDD در زمینه کشف دانش و داده کاوی، سانفرانسیسکو، کالیفرنیا ۱۳ تا ۱۷ آگوست ۲۰۱۶.

۷. استرومیل، آی. و کونونکو، آی. توضیح مدل های پیش بینی و پیش بینی های منفرد با مشارکت ویژگی ها. *دانش. سیستم اطلاعاتی* ۶۴۷، ۶۶۵-۲۰۱۴.

۸. لیپتون، زد. سی. افسانه تفسیرپذیری مدل. *ص. ACM* ۱۶، ۳۱-۵۷-۲۰۱۸.

۹. هاستی، ت.، تیپشیرانی، ر. و فریدمن، ج. عناصر یادگیری آماری: داده کاوی، استنتاج و پیش بینی (انتشارات اشپرینگر نیویورک، نیویورک، ۲۰۰۸).

۱۰. جرز، جی. ام. و همکاران جایگذاری داده های گمشده با استفاده از روش های آماری و یادگیری ماشین در یک مسئله واقعی سرطان سینه. *مصنوع. هوشمند. پزشکی* ۱۵، ۵۰-۱۱۵-۲۰۱۰.

و جامعه فارماکومتری. مقدمه مختصر ما با طیف وسیعی از منابع مرتبط تکمیل شده است. ما با ذکر مثال هایی مرتبط با توسعه دارو، زمینه را فراهم کرده ایم. در پایان، خلاصه ای از وضعیت فعلی حوزه های یادگیری ماشینی و فارماکولوژی بالینی و همچنین چشم اندازی از چگونگی ادغام بیشتر این حوزه ها در آینده ارائه می دهیم. روش های آماری پیشرفته برای فارماکومتری دانان جدید نیستند؛ در واقع، چنین روش هایی مدتی است که برای توصیف پدیده های PK و PD استفاده می شوند. به عنوان مثال، روش های بیزی یک جزء جافاده از رویکردهای فارماکومتری هستند. ۷۲، ۷۳.

بنابراین، به نظر می رسد که با تثبیت و برجسته تر شدن رویکردهای آماری و یادگیری ماشینی در صنعت داروسازی، متخصصان فارماکومتری از جمله کسانی خواهند بود که از این روش ها بهره می برند. علاوه بر این، فرصت های جدیدی برای بررسی سایر سؤالات بالینی، مانند طبقه بندی بیمار بر اساس ویژگی های پایه با ابعاد بالا، ممکن است در فارماکولوژی بالینی با استفاده از رویکردهای یادگیری ماشینی امکان پذیر شود.

چندین نمونه از رویکردهای یادگیری ماشینی که در سؤالات فارماکولوژی بالینی به کار گرفته شده اند، شامل ادغام تکنیک های مدل سازی «کلاسیک»، مانند تعیین یک مدل ساختاری مبتنی بر درک مکانیسمی، و رویکردهای یادگیری ماشینی است. ۶۹-۷۱ رویکردهای فارماکومتری کلاسیک مبتنی بر اصول فارماکولوژیکی هستند که فرضیه های حاصل از درک فیزیولوژی و خواص دارو را منعکس می کنند. بعید است که این مدل ها در آینده نزدیک به طور کامل با رویکردهای یادگیری ماشینی جایگزین شوند. با این حال، هنگامی که مجموعه داده ها و مسائل پیچیده تر هستند، تأثیرات و روابط ناشناخته زیادی وجود دارد و تمرکز بر درون یابی و ارزیابی سریع است، فارماکومتری ممکن است از به کارگیری روش هایی از نوع یادگیری ماشینی بهره مند شود. در ادامه، انتظار داریم که تلفیق این درک با مدل های یادگیری ماشینی بتواند به مدل های بسیار مؤثری در آینده منجر شود. یک مقاله اخیر با دیدگاهی جدید، جزئیات بیشتری در مورد کاربردهای یادگیری ماشینی در فارماکولوژی بالینی ارائه می دهد. ۷۴ در عصر کلان داده، فرصت های جدید زیادی برای یادگیری ماشینی در داروشناسی بالینی وجود دارد. به عنوان مثال، داده های تولید شده از دستگاه های پوشیدنی، چالش های جدیدی را در مورد چگونگی پیوند آنها با داده های PK در آینده ایجاد می کنند. علاوه بر این، دسترسی به داده های دنیای واقعی می تواند شواهد محکمی برای تغییرهای کمی، مجموعه داده های کنترل مکمل و مدل های تقویت کننده ای که روی مجموعه داده های کوچک آموزش دیده اند، فراهم کند.

در رویکردهای فارماکومتری، یک مدل پیش بینی کننده معمولاً با ادغام یک مدل ساختاری و داده های مرتبط ایجاد می شود. مدل ساختاری به طور قابل توجهی فضای راه حل را محدود می کند و بنابراین، داده های نسبتاً کمی برای برازش مدل مورد نیاز است. برعکس، در شبکه های عصبی، ساختار مدل از پیش مشخص نشده است و بنابراین، داده های نسبتاً بیشتری برای ساخت یک مدل پیش بینی کننده مورد نیاز است. همچنین لازم به ذکر است که ما هنوز در مرحله ابتدایی درک هستیم که در آن نقطه ادغام داده های بزرگتر با این روش های جدید یادگیری ماشین می تواند برای عملکرد در مقایسه با روش های سنتی تر مفید باشد. چالش زیره پیش بینی سری های زمانی نشان می دهد که ترکیب روش های آماری کلاسیک و یادگیری ماشینی دقیق ترین پیش بینی را ایجاد می کند و بنابراین، آن را به عنوان راهی برای پیشرفت پیشنهاد می دهد. یکی از محرک های اصلی موفقیت رویکردهای فارماکومتری این است که مدل ها شامل درک کاملی از فرآیندهای جذب، توزیع، متابولیسم و حذف دارو هستند. مدل های ایجاد شده بسیار پیش بینی کننده هستند و بنابراین، کاربرد گسترده ای در پشتیبانی از دارو پیدا می کنند.

۱۱. دورگو، ای. وی، ارشوف، وی. و گولین، ای. کت پوست: تقویت گردانان با پشتیبانی از ویژگی های دسته بندی شده. پیش چاپ *arXiv arXiv:1810.11363* (۲۰۱۸).
- نظریه تصحیح سوگیری انتخاب نمونه. در (Zeugmann, T. Jeds, C., Mohri, M., Riley, M., Rostamizadeh, A., Freund, Y., Györfi, L., Turán, G. & 2. نظریه یادگیری الگوریتمی. 38-53 (اشپرنگر، برلین، هایدلبرگ، 2008).
۱۳. بی. کی، لسلر، جی. و استوارت، ای. ای. بهبود وزن دهی امتیاز گرایش با استفاده از یادگیری ماشین. *آمار، پزشکی*. ۳۷، ۳۴۶-۳۴۰ (۲۰۱۰).
۱۴. نیوبی، دی.، فریتاس، ای. ای. و غفورین، تی. مقابله با مجموعه داده های کلاس نامتعادل در مدل های جذب خوراک. *مجله شیمی، مدل اطلاعات*. ۵۳، ۴۶۱-۴۷۴ (۲۰۱۳).
۱۵. هو، ال. اچ.، هوانگ، ام. دلیو، کی. اس. دلیو، و تسای، سی. اف. تأثیر تابع فاصله بر طبقه بندی k-نزدیک ترین همسایه برای مجموعه داده های پزشکی. *اسپرینگر پلاس* ۴، ۱۳۰-۱۳۶ (۲۰۱۶).
۱۶. پروزان، آی.، وات، اس. و فرتی، وی. ادغام معیارهای شباهت توالی مبتنی بر هم ترازی و بدون هم ترازی برای طبقه بندی توالی های زیستی. *بیوانفورماتیک* ۱۴، ۱۳۹۶-۱۴۰۳ (۲۰۱۵).
۱۷. گولیسون، ای. ای.، بیل، پی.، چانگ، دی. کی. و بیانکین، ای. وی. زیرگروه های مولکولی سرطان پانکراس. *نات. Rev. Gastroenterol. هپاتول*. ۱۶، ۲۰-۲۰۷ (۲۰۱۹).
۱۸. لویو، اس. کوانتیزاسیون کمترین مربعات در PCM: نظریه انتقال اطلاعات *IEEE* ۱۳۷-۱۳۹ (۱۹۸۲).
- خوشه بندی مبتنی بر چگالی. در (Kriegel, H.-P., Kröger, P., Sander, J., Zimek, A. 19. وایی ایتنر دیکس. حداقل دانش داده های بازایی شده. دیسکاوا، ۲۳۱-۲۴۰) (۲۰۱۱).
20. استر، م.، کریگل، اچ. پی.، ساندز، جی. و شو، XA الگوریتم مبتنی بر چگالی برای کشف خوشه ها، یک الگوریتم مبتنی بر چگالی برای کشف خوشه ها در پایگاه های داده مکانی بزرگ با نویز. *مجموعه مقالات دومین کنفرانس بین المللی کشف دانش و داده کاوی*. AAAI Press, Portland, OR, 1996 (226-231).
۲۱. زیمک، آ.، شوبرت، ای. و کریگل، اچ. پی. بررسی تشخیص داده های پرت بدون نظارت در داده های عددی با ابعاد بالا. *داده های آماری، تحلیلی، حداقل*. ۵، ۳۶۳-۳۸۷ (۲۰۱۲).
۲۲. پیرسون، ک. درباره خطوط و صفحاتی که بیشترین برازش را به سیستم های نقاط در فضا دارند. *مجله فلسفی و مجله علمی لندن، ادینبورگ و دوبلین*. ۲، ۵۵۹-۵۷۲ (۱۹۰۱).
۲۳. t-SNE. van der Maaten, L. & Hinton, G. 23. جی. ماخ. *پژوهش های یادگیری*. ۹، ۲۵۷۹-۲۶۰۵ (۲۰۰۸).
۲۴. یخت، ای. و همکاران کاهش ابعاد برای مصورسازی داده های تک سلولی با استفاده از UMAP. *نات. بیوتکنولوژی*. ۳۷، ۳۸۰-۳۸۴ (۲۰۱۹).
۲۵. نگوین، ال اچ و هولمز، اس. ده نکته سریع برای کاهش ابعاد مؤثر. *PLoS Comput. Biol.* 15(10): e1006907 (2019).
۲۶. وانگ، سی.، ماچیراجو، آر. و هوانگ، کی. طبقه بندی بیماران سرطان سینه با استفاده از روش خوشه بندی اجماعی منظم مولکولی. *روش ها*. ۴۰، ۳۱۲-۳۲۰ (۲۰۱۴).
۲۷. کوپر، جی. اس، باینوم، ام ال و سامرز، ای. سی. بینش های اخیر در اپیدمیولوژی بیماری های خودایمنی: بهبود تخمین شیوع و درک خوشه بندی بیماری ها. *جی. خودایمنی*. ۳۳، ۱۹۷-۲۰۷ (۲۰۰۹).
۲۸. والش، دی. و ریپیک، ال. خوشه بندی علائم در سرطان پیشرفته. *بشپتیانی، مراقبت سرطان*. ۱۴، ۸۳۱-۸۳۶ (۲۰۰۶).
۲۹. کنسرسیونم ژنوتیپ-بیان بافت (GTEx) و همکاران آنالیز آزمایشی ژنوتیپ-بیان بافت (GTEx): تنظیم ژن چند بافتی در انسان. *علم*. ۳۴۸، ۶۴۸-۶۶۰ (۲۰۱۵).
30. ژانگ، جی. دی، برنتنيس، ان.، راث، ای. و ایلینگ، ام. داده کاوی شبکه ای از ژن های پاسخ اولیه را به عنوان امضای اجماعی سمیت ناشی از دارو در شرایط آزمایشگاهی (in vitro) و درون بدن (in vivo) آشکار می کند. *فارماکوژنومیک*. ۱۴، ۲۰۸-۲۱۶ (۲۰۱۴).
۳۱. جما، آ. و همکاران خوشه بندی داروهای ضد سرطان در سرطان ریه بر اساس پروفایل های بیان ژن و پایگاه داده حساسیت سرطان. *EBMC*. ۱۷۴، ۳۰۶ (۲۰۰۶).
۳۲. کخ، ام. ای و والدمن، اچ. خوشه بندی شباهت ساختار پروتئین و ساختار محصول طبیعی به عنوان اصول راهنما در کشف دارو. *کشف مواد مخدر، امروز*. ۱۰، ۴۷۱-۴۸۳ (۲۰۰۵).
۳۳. روتلینگر، م. و اشنادر، گ. کاهش ابعاد غیرخطی و نگاشت کتابخانه های ترکیبات برای کشف دارو. *مدل، مجله مول*. ۳۴، ۱۰۸-۱۱۷ (۲۰۱۲).
۳۴. عزت، آ.، وو، م.، لی، ایکس. ال. و کو، سی. کی. پیش بینی برهمکنش دارو-هدف با استفاده از یادگیری جمعی و کاهش ابعاد. *روش ها*. ۸۱، ۸۸-۱۰۷ (۲۰۱۷).
۳۵. نستروف، ی. *سخرانی های در مورد بهینه سازی محدب*. جلد ۱۳۷ (اشپرنگر، ۲۰۱۸).
۳۶. کاولی، جی سی و تالوت، ان ال سی، درباره بیش برازش در انتخاب مدل و سوگیری انتخاب بعدی در ارزیابی عملکرد. *جی. ماخ. یادگیری. پژوهش*. ۱۱، ۲۰۷۹-۲۱۰۷ (۲۰۱۰).
۳۷. میچل، تی ام یادگیری ماشین (شرکت مک گرا-هیل، نیویورک، ۱۹۹۷).
۳۸. بلسون، دلیو. ای. تطبیق و پیش بینی بر اساس اصل طبقه بندی بیولوژیکی. *آمار R/ جامعه آماری*. گروه C. *آمار کاربردی*. ۸، ۶۵-۷۵ (۱۹۵۹).
۳۹. کریشر، جی. پی. کتابشناسی مشروح از کاربردهای تحلیل تصمیم گیری در مراقبت های بهداشتی. *عملیات، تحقیق*. ۲۸، ۹۷-۱۱۳ (۱۹۸۰).
۴۰. شورتلیف، ای اچ، بوکانان، بی جی و فیگنباوم، ای. ای. مهندسی دانش برای تصمیم گیری پزشکی: مروری بر ابزارهای تصمیم گیری بالینی مبتنی بر کامپیوتر. *مجموعه مقالات IEEE* ۷، ۱۲۰۷-۱۲۲۴ (۱۹۷۹).
۴۱. باخ، پی اچ و بریجز، جی دلیو. رویکرد درخت تصمیم گیری برای کاربرد مطالعات متابولیسم سینتیک دارو در آزمایش های سم شناسی و داروشناسی درون تنی و برون تنی. *توکسیکول، مکمل*. ۸، ۱۷۳-۱۸۸ (۱۹۸۵).
۴۲. جردن، تی جی و رایشمن، ال بی. دوزاز یک بار در روز در مقابل دوز دو بار در روز تزئوفیلین: رویکردی برای تحلیل تصمیم گیری جهت ارزیابی سطح خونی تزئوفیلین و میزان پذیرش. *جانب کشیش، تنفس، دیس*. ۱۴، ۱۵۷۳-۱۵۷۷ (۱۹۸۹).
۴۳. بریمن، ل.، فریدمن، ج.، اولشن، ر. و استون، س. درخت های طبقه بندی و رگرسیون. *گروه بین المللی وادزورث*. ۳۷، ۳۲۷-۳۵۱ (۱۹۸۴).
۴۴. کوینلان، جی آر. القای درخت های تصمیم گیری. *ماخ. یاد بگیر*. ۱۱، ۱-۶ (۱۹۸۶).
۴۵. بریمن، ل. جنگل های تصادفی. *ماخ. یاد بگیر*. ۵، ۴۵-۳۲ (۲۰۰۱).
۴۶. سگال، معیارهای یادگیری ماشین MR و رگرسیون جنگل تصادفی. گزارش فنی، (مرکز بیوانفورماتیک و آمار زیستی مولکولی، دانشگاه کالیفرنیا، سانفرانسیسکو، کالیفرنیا، ۲۰۰۳).
47. Ben-Hur, A., Ong, CS, Sonnenburg, S., Schölkopf, B. & Rätsch, J. ماشین های بردار پشتیبان و هسته ها برای زیست شناسی محاسباتی. *PLoS Comput. Biol.* 4(10): e1000173 (2008).
۴۸. نیومن، آ. و همکاران شمارش دقیق زیرمجموعه های سلولی از پروفایل های بیان بافت. *روش های طبیعی*. ۱۲، ۴۵۳-۴۵۷ (۲۰۱۵).
۴۹. واپنیک، وی. و لرنر، ای. تشخیص الگو با استفاده از روش پرتله تعمیم یافته. *کنترل اتوماتیک، رم*. ۲۴، ۷۷۴-۷۸۰ (۱۹۶۳).
۵۰. اسمولا، ای جی و شولکوف، بی. آموزشی در مورد رگرسیون بردار پشتیبان. *آمار، محاسبات*. ۱۴، ۱۹۹-۲۲۲ (۲۰۰۴).
۵۱. مرسر، جی. توابع از نوع مثبت و منفی و ارتباط آنها با نظریه معادلات انتگرال. *فلسفه، ترجمه، آر. جامعه شناسی، لندن، بی. زیست، علمی*. ۴۱۵-۴۴۶ (۱۹۰۹).
52. Schölkopf, B. & Smola, AJ. *یادگیری با هسته ها: ماشین های بردار پشتیبان، منظم سازی، بهینه سازی و فراتر از آن* (انتشارات MIT، کمبریج، ماساچوست، ۲۰۰۲).
۵۳. کریزفسکی، آ.، ساتسکور، آ. و هینتون، جی. ای. طبقه بندی ایمپج نت با شبکه های عصبی کانولوشن عمیق. *پیشرفت ها در سیستم های پردازش اطلاعات عصبی*. ۲۵، ۱۰۹۷-۱۱۰۵ (۲۰۱۲).
۵۴. چو، ک. و همکاران یادگیری نمایش عبارات با استفاده از رمزگذار-رمزگشای RNN برای ترجمه ماشینی آمار. *پیش چاپ arXiv arXiv:1406.1078* (۲۰۱۴).
۵۵. چو، ک.، ون مریبوئر، ب.، باهداناو، د. و بنجیو، ی. درباره ی ویژگی های ترجمه ی ماشینی عصبی: رویکردهای رمزگذار-رمزگشا. *arXiv:1409.1259* (۲۰۱۴).
۵۶. کینگما، دی پی، محمد، اس، رزنده، دی جی و ولینگ، ام. یادگیری نیمه نظارتی با مدل های مولد عمیق. NIPS'14: مجموعه مقالات بیست و هفتمین کنفرانس بین المللی سیستم های پردازش اطلاعات عصبی، ۳۵۸۱-۳۵۸۹ (۲۰۱۴).
۵۷. چوی، ای. و همکاران استفاده از مدل های شبکه عصبی بازگشتی برای تشخیص زودهنگام شروع نارسایی قلبی. *دانشگاه OUP* (24/2/361/2631499) <https://academic.oup.com/jamia/article> (2019).
۵۸. تانگ، جی. تی.، کائو، وای، شیائو، جی. وای. و گوئو، کیو. ال. پیش بینی غلطت پلاسمای ریفنتانیل بر اساس شبکه عصبی ال. *دانشگاه جی. سنت جنوبی*. ۳۱۸۷، ۳۱۹۲-۳۱۹۲ (۲۰۱۳).
59. پاسکانو، ر.، میکولوف، ت. و بنجیو، ی. در مورد دشواری آموزش شبکه های عصبی بازگشتی. کنفرانس بین المللی یادگیری ماشین، ادینبورگ، 26 ژوئن - 1 ژوئیه 2013 (2013).

۶۹. بنهامو، پ.ی. و همکاران: تحویل انسولین با حلقه بسته در بزرگسالان مبتلا به دیابت نوع 1 در شرایط واقعی: یک کارآزمایی متقاطع تصادفی کنترل شده چند مرکزی 12 هفته ای با برچسب باز/نست دیگ سلامت. (2019) e17-e25
70. چاوادا، آر.، قوش، ن.، ساندارادورا، آی.، مالی، ام. و ون هال، تاسیس یک AUC توسط JZ-4: آستانه نفروتنوکسیسییتی گامی به سوی دوز و انکومایسین فردی برای باکتریی استافیلوکوکوس اورئوس مقاوم به متی سلین است. ضد میکروب. عوامل شیمی درمانی. (2017) e02535-16
۷۱. هنینگ، س.، هولتهوس، ف. و استاتز، س. ای. مقایسه روش های تنظیم دوز برای توپرامایسین یک بار در روز در بیماران کودک و نوجوان مبتلا به فیبروز کیستیک. کلینیک. فارماکوکینتیک. ۵۴، ۴۲۱-۴۰۹ (۲۰۱۵).
۷۲. دانسیریکول، سی.، موریس، آر. جی.، تت، اس. ای. و دافول، اس. بی. رویکرد بیزی برای مدل سازی فارماکوکینتیک جمعیتی سیرولیموس. *مجله پزشکی بالینی بریتانیا (Br. J. Clin. Pharmacol.)* ۴۲، ۴۳۴-۴۲۰ (۲۰۰۶).
۷۳. لان، دی جی، پست، ان.، توماس، ای.، ویکفیلد، جی. و اشیپگل هالتر، دی. تحلیل بیزی مدل های PK/PD جمعیتی: مفاهیم کلی و نرم افزار. *مجله فارماکوکینتیک. فارماکودین.* ۲۹، ۳۰۷-۳۰۲ (۲۰۰۲).
۷۴. هچینسون، ل. و همکاران: مدل ها و ماشین ها: چگونه یادگیری عمیق، داروشناسی بالینی را به سطح بعدی خواهد برد. *Pharmacomet. System. Pharmacol. CPT* ۱۳۱، ۱۳۴-۱۳۱ (۲۰۱۹).
۷۵. ماکریداکیس، س.، اسپیلیوتیس، ای. و آسیماکوپولوس، و. رقابت M4: نتایج، یافته ها، نتیجه گیری و مسیر پیش رو. *پیش بینی بین المللی*، ۸۰۸-۸۰۲، ۳۴ (۲۰۱۸).
- در زمینه آکوستیک، گفتار و پردازش سیگنال، کیوتو، 25-30 مارس، 8624-8628 IEEE پیشرفت ها در بهینه سازی شبکه های بازگشتی. کنفرانس بین المللی. Bengio, Y., Boulanger-Lewandowski, N. & Pascanu, R. (2012)
۶۱. هوخرایتر، س. و اشمیدهور، ج. حافظه کوتاه مدت بلندمدت. *محاسبات عصبی*، ۹، ۱۷۳۵-۱۷۸۰ (۱۹۹۷).
۶۲. گودفلو، آی.، بنگیو، وای. و کورویل، ای. یادگیری عمیق (انتشارات MIT، کمبریج، ۲۰۱۶) <http://www.deeplearningbook.org> (MA <2016>)
63. چانگ، جی.، گولسهره، سی.، چو، کی. و بنگیو، وای. ارزیابی تجربی شبکه های عصبی بازگشتی دروازه دار در مدل سازی توالی. کارگاه آموزشی یادگیری عمیق در NIPS، مونترال، 8 تا 13 دسامبر (2014).
۶۴. ریو، جی وای، کیم، اچ یو و لی، اس وای. یادگیری عمیق پیش بینی تداخلات دارو و غذا را بهبود می بخشد. *مجموعه مقالات علمی، آکادمی ملی علوم* ۱۱۵، ۲۰۱۸ (E4304-E4311)
65. Hammann, F., Gutmann, H., Vogt, N., Helma, C. و Drewe, جی. پیش بینی عوارض جانبی داروها با استفاده از مدل سازی درخت تصمیم گیری. *کلینیک. فارماکول. درمان*، ۵۲، ۵۸-۵۹ (۲۰۱۰).
66. لنکستر، ام سی و سوبی، ای. پیش بینی بهبود یافته ی تورساد دو پوینت های ناشی از دارو از طریق شبیه سازی دینامیک و الگوریتم های یادگیری ماشین. *کلینیک. فارماکول. درمان*، ۳۷۱، ۱۰۰-۳۷۹ (۲۰۱۶).
۶۷. دانهاور، آی. و همکاران: پیش بینی زود هنگام بهبود یافته هیپر بیلیروبینمی نوزادی مرتبط با علائم بالینی با یادگیری ماشین. *اطفال. پژوهش*، ۱۲۲، ۸۶ (۲۰۱۹).
۶۸. گاودا، ای. و همکاران: شخصی سازی مدیریت دارویی کم خونی با استفاده از یادگیری تقویتی. *شبکه های عصبی*، ۱۸، ۸۲۶-۸۳۴ (۲۰۰۵).