

کاربرد روش ماشین‌های بردار پشتیبان در طبقه‌بندی داده‌های اقلیمی

سیدمسعود سلیمان‌پور^۱، پارسا حقیقی^۲

۱. دانشیار بخش تحقیقات حفاظت خاک و آبخیزداری، مرکز تحقیقات و آموزش کشاورزی و منابع طبیعی استان فارس، سازمان تحقیقات، آموزش و ترویج کشاورزی، شیراز، ایران
(نویسنده مسؤل). رایانامه: m.soleimanpour@areeo.ac.ir

۲. کارشناس بخش تحقیقات حفاظت خاک و آبخیزداری، مرکز تحقیقات و آموزش کشاورزی و منابع طبیعی استان فارس، سازمان تحقیقات، آموزش و ترویج کشاورزی، شیراز، ایران. رایانامه: p.haghighi@areeo.ac.ir

تاریخ دریافت: ۱۴۰۲/۰۸/۰۳

تاریخ ویرایش: ۱۴۰۲/۰۹/۲۵

تاریخ پذیرش: ۱۴۰۲/۱۱/۱۰

تاریخ چاپ: ۱۴۰۳/۰۱/۲۰

صص: ۱۱-۱

چکیده

در یک ساختار اقلیمی، شناسایی، پیش‌بینی و مدیریت بحران اهمیت زیادی دارد. برای درک نظام‌های پیچیده و شبیه‌سازی و پیش‌بینی رفتار آن‌ها از الگوها، به‌عنوان ابزارهایی کاربردی، استفاده می‌شود. ماشین‌های بردار پشتیبان یکی از روش‌های یادگیری نظارت‌شده است که برای طبقه‌بندی و رگرسیون از آن استفاده می‌شود. ماشین‌های بردار پشتیبان می‌توانند الگوهای پنهان را تشخیص داده و به تغییرات پیچیده در داده‌های اقلیمی پاسخ دهند. در این مقاله، ساختار روش ماشین‌های بردار پشتیبان و کاربرد آن‌ها در طبقه‌بندی داده‌های اقلیمی بررسی شده است. ویژگی‌های ساختار ماشین‌های بردار پشتیبان به انتخاب نوع تابع هسته مرتبط است. بنابراین، در انتخاب نوع تابع هسته باید دقت زیادی شود. از سوی دیگر، در پیش‌بینی اقلیمی، تعیین مؤلفه اصلی مرحله بسیار مهمی است تا با استفاده از تعداد مؤلفه بهینه بتوان بین داده‌های پیش‌بینی‌کننده و پیش‌بینی‌شونده بهترین برآزش را انجام داد.

کلیدواژه‌ها: ماشین‌های بردار پشتیبان، طبقه‌بندی داده‌های اقلیمی، توابع هسته، تعیین مؤلفه اصلی.

مقدمه

ایران با مشکلات بی‌سابقه اقلیمی از جمله خشک شدن دریاچه‌ها و رودخانه‌ها، توفان‌های گردوغبار، دمای بی‌سابقه، خشک‌سالی و سیل روبرو است. نگرانی از تغییرات آب‌وهوایی و تاثیرات آن بر جنبه‌های مختلف زندگی انسان در حال افزایش است. در مدیریت خشک‌سالی و سیلاب، پیش‌بینی‌های اقلیمی بسیار بااهمیت است. بنابراین، لزوم توجه به روش‌های نوین در پیش‌بینی اقلیمی و بررسی کارکرد روش‌ها بر کسی پوشیده نیست. پیش‌بینی فرآیندهای آب و هوایی، به‌ویژه پیش‌بینی بارندگی، ابزار مناسبی در اختیار مدیران بخش‌های مختلف (کشاورزی و منابع طبیعی و غیره) قرار می‌دهد تا برای بهینه‌سازی هزینه‌ها، امکانات و بهره‌وری حداکثری از آن‌ها، مقابله با خشک‌سالی و کاهش خسارت ناشی از آن، سیاست‌های آینده را برنامه‌ریزی کنند. بارندگی مبحثی کلیدی در چرخه آشناسی است که بیش‌ازپیش، هم در سطح منطقه و هم در سطح جهان، در حال تغییر است (نسوبوگا و همکاران^۱، ۲۰۱۴) و درزمینه مدیریت منابع آب یکپارچه، پیش‌بینی نیازهای محصول و ارزیابی محیط‌زیست به شکلی گسترده استفاده می‌شود (هیان و همکاران^۲، ۲۰۰۷). بنابراین، برای مدیریت منابع آب، کشاورزی و بسیاری از کاربران، مؤلفه‌های آب و هوایی بسیار جالب توجه است (گابین و همکاران^۳، ۲۰۱۳). همچنین، موقعیت جغرافیایی ایران در منطقه کناره حاره، کشور ما را در کمربند خشک جهان قرار داده است. باتوجه به این موقعیت جغرافیایی، میانگین بارش ایران بسیار کمتر از میانگین جهانی آن است (صفوی، ۱۳۸۸). یادگیری ماشینی^۴ به‌عنوان یکی از شاخه‌های وسیع و پرکاربرد هوش مصنوعی، شیوه‌ها و روش‌هایی را تنظیم و ابداع می‌کند که رایانه‌ها و سامانه‌ها براساس آن‌ها می‌توانند یاد بگیرند. استفاده از یادگیری ماشینی در کاهش هزینه‌های عملیاتی و بهبود سرعت عمل تجزیه و تحلیل داده‌ها بسیار به‌صرفه است.

یادگیری ماشینی به دو دسته کلی «یادگیری نظارت‌شده»^۵ و «یادگیری نظارت‌نشده»^۶ تقسیم می‌شود. در روش‌های یادگیری

ماشینی نظارت‌شده، مجموعه‌ای از بردارهای ورودی مانند X و بردارهای خروجی متناظر با آن‌ها مانند t به رایانه ارائه می‌شود. هدف آن است که با استفاده از این داده‌های آموزشی در ورودی X جدید، ماشین بتواند t را پیش‌بینی کند. از جمله روش‌های یادگیری نظارت‌شده می‌توان به روش‌های طبقه‌بندی^۷ مانند شبکه‌های عصبی^۸، درخت تصمیم^۹، بیز ساده^{۱۰}، نزدیک‌ترین همسایگی^{۱۱}، ماشین‌های بردار پشتیبان^{۱۲} و روش‌های رگرسیون مانند رگرسیون خطی^{۱۳}، رگرسیون غیرخطی^{۱۴} و رگرسیون بردار پشتیبان^{۱۵} اشاره کرد (استوارت و همکاران^{۱۶}، ۲۰۱۰). در روش‌های یادگیری نظارت‌نشده، یادگیری ماشین تنها از طریق داده‌های ورودی انجام می‌شود. یعنی، مجموعه داده‌ها تنها شامل متغیرهای ورودی است و هیچ خروجی متناسبی با ورودی‌ها وجود ندارد. بنابراین، در یادگیری نظارت‌نشده، الگوریتم یادگیری، خود، الگو و ساختار میان داده را پی‌می‌گیرد. در واقع، یادگیری نظارت‌نشده روشی است که از آن برای یافتن الگوهای میان داده‌ها استفاده می‌شود. به‌عبارت‌دیگر، با یادگیری نظارت‌نشده می‌توان ساختار و الگوهای پنهان میان داده‌ها را یافت. از جمله روش‌های یادگیری نظارت‌نشده می‌توان به روش‌های خوشه‌بندی^{۱۷} مانند کی-میانگین^{۱۸}، دی بی اسکن^{۱۹}

1. Nsubuga et al
2. Haiyun et al
3. Guobin et al
4. Machine Learning
5. Supervised Learning
6. Unsupervised Learning
7. Classification
8. Artificial Neural Network (ANN)
9. Decision Tree
10. Naive Bayes
11. K Nearest Neighbor
12. Support Vector Machines (SVM)
13. Linear Regression
14. Non-Linear Regression
15. Support Vector Regression
16. Stuart J at al
17. Clustering
18. K-Medoids
19. DBSCAN

مجموعه‌ای از داده‌هایی که خروجی مطلوب آن‌ها از قبل مشخص است، الگو را آموزش می‌دهند. وقتی یک داده جدید که خروجی آن مشخص نیست به الگو ارائه شد، الگو می‌تواند خروجی مطلوب را تولید کند. ولادیمیر وپنیک، از محققین روسی، در سال ۱۹۶۵ نظریه آماری یادگیری را به صورت مستحکم‌تری بنا نهاد و بر این اساس ماشین‌های بردار پشتیبان را ارائه داد. ماشین‌های بردار پشتیبان دارای ویژگی‌های زیر هستند (شین و همکاران^۸، ۲۰۰۵):

- طراحی طبقه‌بند با حداکثر تعمیم؛
- رسیدن به نقطه بهینه کلی تابع؛
- تعیین خودکار ساختار و جاینگاری بهینه برای طبقه‌بندی‌کننده؛
- الگو کردن توابع تمایز غیرخطی با استفاده از هسته‌های غیرخطی و مفهوم حاصل ضرب داخلی در فضاها هیلبرت (شین و همکاران، ۲۰۰۵).

ماشین‌های بردار پشتیبان الگوریتمی است که نوع خاصی از الگوهای خطی را می‌یابد که حداکثر حاشیه ابرصفحه را به دست می‌دهند. حداکثر کردن حاشیه ابرصفحه به حداکثر شدن تفکیک بین طبقات منجر می‌شود. به نزدیک‌ترین نقاط آموزشی به حداکثر حاشیه ابرصفحه، بردارهای پشتیبان اطلاق می‌شود. اگر داده‌ها به صورت خطی از یکدیگر مجزا باشند تنها از این بردارها (نقاط) برای مشخص کردن مرز بین طبقات استفاده می‌شود. ماشین‌های بردار پشتیبان به ماشین‌های خطی آموزش می‌دهند که چگونه سطح بهینه‌ای تولید کنند که بدون خطا و با حداکثر فاصله میان صفحه و نزدیک‌ترین نقاط آموزشی بردارهای (پشتیبان)، داده‌ها را تفکیک کنند (مین

و روش‌های کاهش ابعاد^۱ مانند تحلیل مؤلفه اصلی^۲ و آنالیز تشخیصی خطی^۳ اشاره کرد (هیستی در همکاران^۴، ۲۰۰۱).

روش ماشین‌های بردار پشتیبان، ابزاری جدید و قدرتمند برای تدوین راه‌حل‌هایی برای طبقه‌بندی و رگرسیون است. این روش بر تئوری یادگیری آماری استوار است و در آن، بر مبنای مجموعه‌ای از داده‌های آموزشی، یادگیری انجام می‌شود. تحقیقات نشان می‌دهند که استفاده از این روش یادگیری، در مقایسه با سایر روش‌های داده‌محور پیشین از جمله شبکه‌های عصبی مصنوعی، حتی در شرایطی که ابعاد داده‌ها زیاد و تعداد نمونه‌ها کم باشد، دقت الگو، نتایج را بهبود می‌بخشد. به طور کلی، از ماشین‌های بردار پشتیبان، هم در دسته‌بندی و هم در رگرسیون استفاده می‌شود. گرچه از مبنای تئوری این روش برای طبقه‌بندی استفاده می‌شود، اما اخیراً با گسترش روش ماشین‌های بردار پشتیبان و مقبولیت آن در استفاده، نسبت به روش‌های قدیمی‌تر از جمله شبکه‌های عصبی مصنوعی، تئوری جدید و جامع تری در محدوده مسائل رگرسیون نیز گسترش یافته است. در سال ۱۹۹۵، وپنیک^۵ استفاده از ماشین‌های بردار پشتیبان، به عنوان ابزار رگرسیون، را برای اولین بار مطرح کرد. از روش ماشین‌های بردار پشتیبان در پیش‌بینی و طبقه‌بندی شاخص‌های اقلیمی استفاده شده است. خاشعی سیوکی و همکاران (۱۳۹۷)، سهندرا^۶ و همکاران (۲۰۱۳) و لو و کین^۷ (۲۰۱۴) کارایی روش ماشین‌های بردار پشتیبان را در طبقه‌بندی و پیش‌بینی مؤلفه‌های بارش و دما بررسی کرده و آن را با سایر روش‌ها مقایسه کرده‌اند. نتایج پژوهش آن‌ها گویای دقت بالای روش ماشین‌های بردار پشتیبان در طبقه‌بندی و پیش‌بینی اقلیمی است.

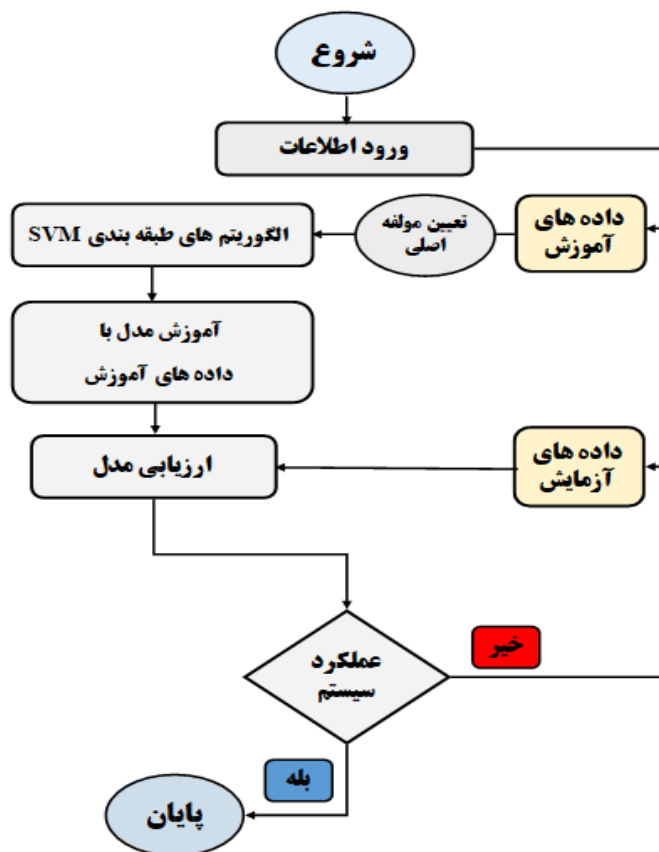
ماشین‌های بردار پشتیبان

ماشین‌های بردار پشتیبان یکی از روش‌های یادگیری با نظارت است که از آن برای طبقه‌بندی و رگرسیون استفاده می‌شود. به طور کلی، در یادگیری نظارت‌شده، با استفاده از

1. Dimensionality Reduction
 2. Principal Component Analysis
 3. Linear Discriminant Analysis
 4. Hastie et al
 5. Vapnik
 6. Sahindra
 7. Lu & Qin
 8. Shin et al
 9. Min et al

ورودی و نوع الگوریتم طبقه‌بندی ماشین‌های بردار پشتیبان بازبینی شده و اجرای الگو تکرار شود. در بحث پیش‌بینی اقلیمی نیز می‌توان از ماشین‌های بردار پشتیبان استفاده کرد. پس از تعیین مؤلفه‌های اصلی، داده‌های مؤلفه‌های اقلیمی به دو دسته آموزش و آزمایش تقسیم می‌شوند. پس از انتخاب نوع تابع هسته و الگوریتم‌های طبقه‌بندی ماشین‌های بردار پشتیبان، خروجی الگو ارزیابی خواهد شد. در شکل ۱ روند کلی الگوسازی ماشین‌های بردار پشتیبان ارائه شده است.

و همکاران^۹، ۲۰۰۵). برای الگوسازی در ماشین‌های بردار پشتیبان، ابتدا داده‌ها به دو بخش آموزش و آزمایش تقسیم شده و اطلاعات پیش‌پردازش می‌شوند. سپس، با استفاده از معیارهای خطا، نوع الگوریتم مناسب برای طبقه‌بندی اطلاعات انتخاب و الگو ارزیابی می‌شود. در ادامه، برای بررسی کارایی الگو، داده‌های آزمایشی به‌عنوان ورودی الگو در نظر گرفته می‌شوند. اگر دقت الگو کافی باشد روند خروجی و گزارش نهایی تکمیل می‌شود، در غیر این صورت باید مؤلفه‌های



شکل ۱. فرآیند الگوسازی در ماشین‌های بردار پشتیبان

ساختار ماشین‌های بردار پشتیبان

ماشین بردار پشتیبان، یک الگوریتم یادگیری با ناظر است که نمونه داده‌ها را به صورت نقاطی در فضا نشان می‌دهد و با استفاده از یک خط^۱ از هم جدا می‌کند. این جداسازی به گونه‌ای است که نقاط داده‌ای که در یک طرف خط هستند مشابه به هم و در یک گروه قرار می‌گیرند. نمونه داده‌های

جدید نیز بعد از اضافه شدن به همان فضا، در یکی از دسته‌های موجود قرار خواهند گرفت. فاصله بین دسته‌های مختلف داده‌ها را حاشیه^۲ گویند و برای تعیین محل مرز بین دو دسته‌بندی، از

1. Hyperplane
2. Margin

پشتیبان به صورت رابطه ۱ است:

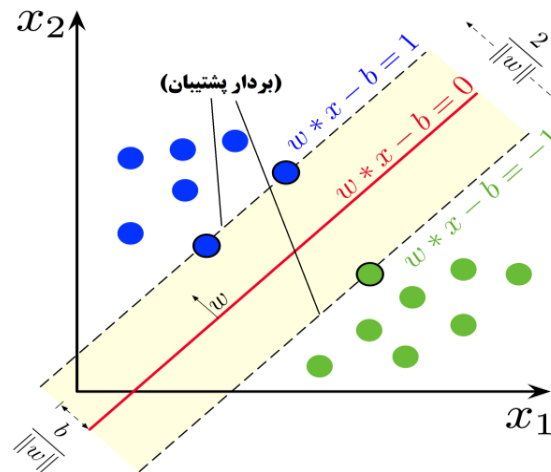
$$w^T * x - b = 0 \quad (1)$$

که در آن، w بردار وزن است که برابر با جمع ضرایب نمونه‌های پشتیبان، x بردار ویژگی‌ها و b عبارت است از تعداد منفی ضرایب نمونه‌های پشتیبان تقسیم بر تعداد نمونه‌های پشتیبان. بردار پشتیبان در ماشین‌های بردار پشتیبان، نقاط حاشیه دسته‌بندی را تعیین می‌کند و به‌عنوان نقاط حاشیه اصلی در ساخت مکان مرز بین دسته‌بندها شناخته می‌شود. در شکل ۲ حاشیه بین دودسته و بردار پشتیبان نشان داده شده است:

بردار پشتیبان استفاده می‌شود. توابع هسته، با اضافه کردن ابعاد بیشتر به مسئله، مسائل تفکیک‌ناپذیر را به مسائل تفکیک‌پذیر تبدیل می‌کند تا دسته‌بندی بهتری حاصل آید. در ادامه ساختار ماشین‌های بردار پشتیبان شرح داده شده است.

۱) بردارهای پشتیبان

بردار پشتیبان در ماشین‌های بردار پشتیبان نقاط داده‌هایی است که برای تصمیم‌گیری در فرآیند طبقه‌بندی استفاده می‌شود. این بردارها به‌طور خاص در مرز جداکننده بین دو کلاس قرار دارند و در تعیین حاشیه جداکننده نقش اصلی را برعهده دارند. فرمول رسم خط جداکننده در ماشین‌های بردار



شکل ۲. حداکثر حاشیه بین دو دسته و بردار پشتیبان

w بردار وزن و $\|w\|$ نرم آن است.

$$\text{Margin} = 2 / \|w\| \quad (2)$$

۳) توابع هسته

همان‌طور که قبلاً اشاره شد، ماشین‌های بردار پشتیبان ممکن است برای تقسیم‌بندی داده‌های غیرخطی استفاده شود، درحالی‌که برای داده‌های خطی طراحی شده است. برای اینکه طبقه‌بندی بر روی داده‌های غیرخطی کار کند

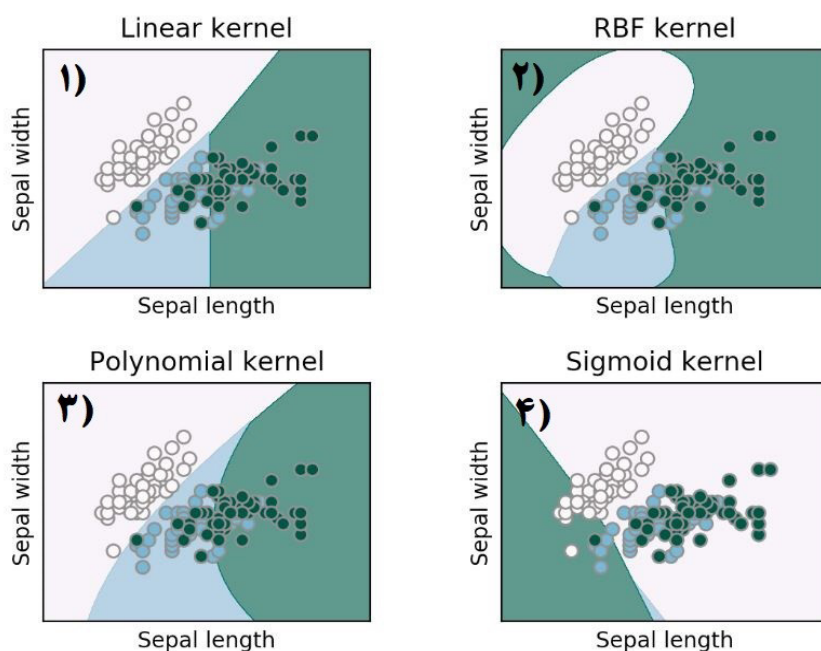
۲) حاشیه

در ماشین‌های بردار پشتیبان، به فاصله بین دسته‌های مختلف داده‌ها در فضای ویژگی، حاشیه می‌گویند. حاشیه بزرگ‌تر نشان‌دهنده جداپذیری بهتر بین دسته‌ها است. در ماشین‌های بردار پشتیبان، هدف اصلی آن است که بین داده‌های دو کلاس مختلف، یک خط جداکننده بیابند. این خط جداکننده باید طوری باشد که فاصله بین خط و نزدیک‌ترین نقطه از هر دو کلاس (که به آن‌ها بردار پشتیبان یا support vectors می‌گویند) حداکثر شود (وین و همکاران^۴، ۲۰۱۵). رابطه ۲ مربوط به حاشیه در ماشین‌های بردار پشتیبان است؛ که در آن،

1. Support Vectors
2. Margin
3. Hyper Plane
4. Wen, et al
5. Kernel Functions

ماشین‌های بردار پشتیبان می‌توان به تابع خطی^۱، چندجمله‌ای^۲، تابع پایه شعاعی^۳، گوسین^۴، سیگموئید^۵، تابع لاپلاس^۶ و هسته شبکه عصبی^۷ اشاره کرد (میر عربی و همکاران، ۲۰۱۹). در شکل ۳ انواع طبقه‌بندی مجموعه داده‌ها با استفاده از کرنل‌های مختلف ارائه شده است.

از مؤلفه‌ای به نام توابع هسته استفاده می‌شود. در ماشین‌های بردار پشتیبان، از توابع هسته برای انتقال داده‌ها به فضای بالاتر استفاده می‌شود. در این فضا، دسته‌بندی بهتری انجام می‌شود. تابع هسته غیرخطی عمل می‌کند و امکان می‌دهد تا داده‌های ورودی به فضای با بعد بالاتر منتقل شوند، به طوری که در این فضا، داده‌های جدید قابل تفکیک باشند. از انواع توابع هسته در



شکل ۳. طبقه‌بندی مجموعه داده‌ها با استفاده از کرنل‌های: (۱) خطی، (۲) تابع پایه شعاعی، (۳) چندجمله‌ای و (۴) سیگموئید

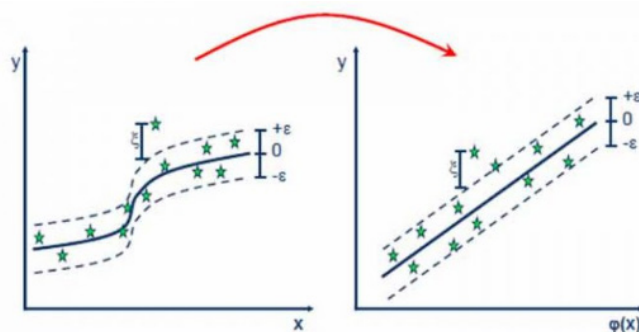
هسته آن بستگی دارد. در ابتدا رگرسیون بردار پشتیبان خطی ارائه شده بود که بعدها برای مسائل غیرخطی تعمیم داده شد. واپنیک برای طبقه‌بندی داده‌های غیرخطی از نگاهت غیرخطی کمک گرفت. برای تعمیم رگرسیون بردار پشتیبان به مسائل غیرخطی، ابتدا داده‌ها با کمک یک نگاهت غیرخطی به فضای

رگرسیون بردار پشتیبان^۸

شکل پر کاربرد ماشین‌های بردار پشتیبان، رگرسیون بردار پشتیبان است. ماشین‌های بردار پشتیبان برای انواع الگوهای پیش‌بینی در طبقه‌بندی و رگرسیون استفاده می‌شود. زمانی که از رگرسیون استفاده می‌شود به آن رگرسیون بردار پشتیبان می‌گویند. در الگوی مسئله رگرسیونی، ورودی‌ها (مقادیر بارش، دما و تبخیر و تعرق) به یک فضای چندبعدی ترسیم می‌شوند و رگرسیون بردار پشتیبان از طریق اصل به حداقل رساندن ریسک ساختاری^۹ آموزش داده می‌شود. رگرسیون بردار پشتیبان روشی ناپارامتریک است که واپنیک^{۱۰} (۱۹۹۵) توسعه داده است. عملکرد رگرسیون بردار پشتیبان به شدت به

1. Linear Kernel
2. Polynomial Kernel
3. Radial Basis Function
4. Gaussian Kernel
5. Sigmoid Kernel
6. Laplacian Kernel
7. Neural Network Kernel
8. Support Vector Regression (SVR)
9. Structural Risk Minimization (SRM)
10. Vapnik

فضای ابعاد بالاتر با استفاده از یک تابع هسته انجام می‌شود تا در آن فضا خط جداساز بهتر و دقیق‌تر عمل کند. شکل ۴ انتقال داده‌ها از یک ساختار غیرخطی به فضای خطی (تریپاتی^۱ و همکاران، ۲۰۰۶).



شکل ۴. انتقال داده‌ها از یک ساختار غیرخطی به فضای خطی

از این روش، مؤلفه‌های اصلی را انتخاب و در الگوسازی از تعداد مؤلفه‌های بهینه استفاده کرد. برای پیش‌بینی اقلیمی بارش و دما می‌توان از ۲۶ متغیر پیش‌بینی کننده بزرگ مقیاس با نام داده‌های بزرگ مقیاس استفاده کرد. داده‌های NCEP برای نمونه شامل میانگین فشار سطح دریا، رطوبت نسبی در ارتفاع ۸۵۰ هکتوپاسکال، میانگین دما در ارتفاع ۲ متری، رطوبت نسبی در ارتفاع ۵۰۰ هکتوپاسکال و غیره است. با تعیین ارتباط بین داده‌های مشاهداتی ایستگاه محلی (بارش و دما) و داده‌های بزرگ مقیاس NCEP می‌توان در شرایط مختلف، میزان بارش و دما را برای آینده پیش‌نمایی کرد.

برای اجرای الگوهای ماشین‌های بردار پشتیبان می‌توان از ابزارهای متفاوتی مثل کتابخانه Scikit Learn و بسته e۱۰۷۱ استفاده کرد. کتابخانه Scikit Learn متن‌باز، مفید، پرکاربرد و قدرتمند در زبان برنامه‌نویسی پایتون است که برای اهداف یادگیری ماشین استفاده می‌شود. این کتابخانه برای یادگیری ماشین و الگوسازی آماری داده‌ها همچون طبقه‌بندی، رگرسیون،

خطی نگاشت می‌یابند. سپس در فضای جدید رابطه خطی بین داده‌ها در فضای ویژگی و خروجی به دست می‌آید که این رابطه خطی در فضای ویژگی، معادل رابطه غیرخطی بین ورودی و خروجی در فضای اصلی است. نگاشت ورودی‌ها به

انتخاب متغیرهای پیش‌بینی کننده غالب

کارن پیرسون (۱۹۰۱) برای اولین بار فن تجزیه به مؤلفه‌های اصلی را شرح داده است. در روش تحلیلی تحلیل مؤلفه اصلی سعی می‌شود، تعداد مؤلفه‌ها تا حد امکان کاهش یابد و در چند مؤلفه اصلی و اثرگذار خلاصه شود. در تحلیل عاملی بعد از به دست آوردن ماتریس همبستگی، ابتدا باید مشخص کنیم که برای تحلیل عاملی از کدام الگو استفاده کنیم. روش تحلیل مؤلفه‌های اصلی دو کارکرد مهم دارد. اولین کارکرد این روش آن است که عامل‌ها را به صورت مستقیم و بدون برآورد اشتراکات، از ماتریس همبستگی تعیین می‌کند. در این روش برای تبیین حداکثر مقدار واریانس متغیرها، ترکیب خطی آن‌ها برآورد می‌شود. بدین صورت که اولین مؤلفه، بیشترین واریانس متغیرها را تبیین می‌کند. سپس مؤلفه دوم بیشترین مقدار واریانس باقی‌مانده در متغیرها را بعد از مؤلفه اول توضیح می‌دهد و به همین ترتیب تا آخر ادامه می‌یابد. کارکرد دیگر تحلیل مؤلفه‌های اصلی این است که مجموعه‌ای از متغیرهای سنجیده شده را به مجموعه‌ای از ترکیب خطی با حداکثر مقدار واریانس تبدیل می‌کند (گوچیچ و همکاران^۲، ۲۰۱۶). در پیش‌بینی متغیرهای اقلیمی نیز می‌توان با استفاده

1. Tripathi

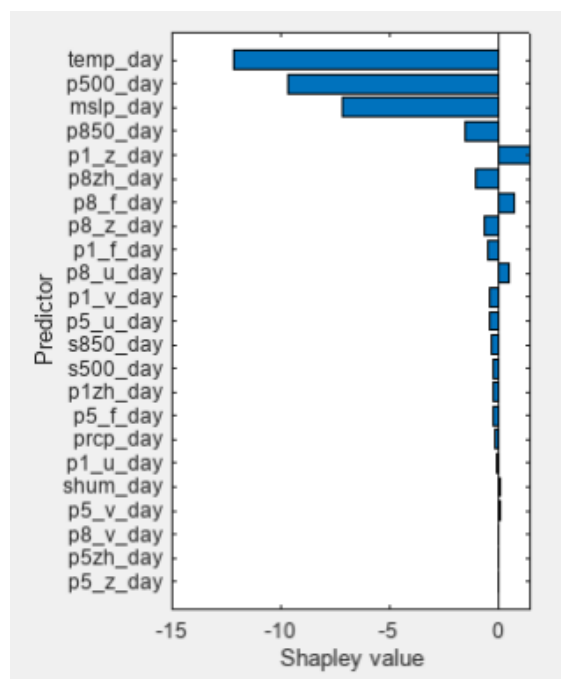
2. Gocic, et al

و ۳۰ درصد در دسته آزمایش قرار گرفتند. برای طبقه‌بندی دمای متوسط از ۲۳ متغیر پیش‌بینی‌کننده بزرگ‌مقیاس با عنوان داده‌های بزرگ‌مقیاس استفاده شد. داده‌های NCEP برای نمونه شامل میانگین فشار سطح دریا، رطوبت نسبی در ارتفاع ۸۵۰ هکتوپاسکال، میانگین دما در ارتفاع ۲ متری، رطوبت نسبی در ارتفاع ۵۰۰ هکتوپاسکال و غیره است. با تعیین ارتباط بین داده‌های مشاهداتی ایستگاه محلی (دمای متوسط ایستگاه سینوپتیک شیراز) و داده‌های بزرگ‌مقیاس NCEP می‌توان در شرایط مختلف میزان دما را برای آینده پیش‌نمایی کرد. برای انتخاب متغیرهای پیش‌بینی‌کننده غالب، با استفاده از تحلیل مؤلفه اصلی، از بین ۲۳ مؤلفه بزرگ‌مقیاس NCEP ۵ مؤلفه انتخاب شدند. شکل ۵ تحلیل مؤلفه اصلی مربوط به ۲۳ مؤلفه بزرگ‌مقیاس NCEP را نشان می‌دهد. بیشترین میزان اثرگذاری مربوط به ۵ مؤلفه: دمای سطح در ارتفاع ۲ متری (temp_day)، میانگین فشار سطح دریا (mslp_day)، ارتفاع ژئوپتانسیلی تراز ۸۵۰ hPa، ارتفاع ژئوپتانسیلی تراز ۵۰۰ hPa و تاوایی سطح (p850_day) است. (p1_z_day)

خوشه‌بندی و کاهش ابعاد، ابزارهای کاربردی زیادی را فراهم می‌آورد. سادگی، قابلیت استفاده، پوشش گسترده الگوریتمی، ارزیابی الگو، تنظیم مؤلفه، پشتیبانی از تقسیم‌بندی داده، پشتیبانی از تعداد زیادی مسئله و غیره از جمله ویژگی‌های مهم این کتابخانه است. بسته ۱۰۷۱e در زبان برنامه‌نویسی R نیز برای الگوسازی ماشین بردار پشتیبان استفاده می‌شود. در این بسته توابع کمکی^۱ نیز وجود دارد. در ساخت ماشین بردار پشتیبان در زبان R و پایتون از رویکرد مشابهی استفاده شده است.

یافته‌ها

در این مطالعه برای طبقه‌بندی داده‌های اقلیمی با روش ماشین‌های بردار پشتیبان از داده‌های ماهانه دمای متوسط ایستگاه سینوپتیک شیراز از سال ۱۹۶۱ تا ۲۰۱۷ استفاده شد و از خروجی الگو CanESM5 در قالب داده‌های NCEP استفاده شده است. داده‌های NCEP شامل ۲۳ مؤلفه بزرگ‌مقیاس اقلیمی است که به صورت روزانه از سال ۱۹۶۱ تا ۲۰۱۷ مربوط به منطقه مورد مطالعه در دسترس است. داده‌ها در دو دسته آموزش و آزمایش تفکیک شدند. ۷۰ درصد داده‌ها در دسته آموزش



شکل ۵. تحلیل مؤلفه اصلی ۲۳ مؤلفه داده‌های بزرگ‌مقیاس NCEP

و آزمایش تحت ۴ نوع توابع هسته متفاوت با هم مقایسه شدند. نتایج ارائه‌شده در جدول ۱ نشان می‌دهد که انواع توابع هسته در طبقه‌بندی داده‌های دمای متوسط ایستگاه سینوپتیک شیراز دقت قابل‌قبولی دارند.

ماشین‌های بردار پشتیبان برای تقسیم‌بندی داده‌ها از انواع مختلف توابع هسته استفاده می‌کنند. در این پژوهش از چهار نوع توابع خطی، تابع پایه شعاعی، چندجمله‌ای و سیگموئید برای طبقه‌بندی اطلاعات استفاده شد. داده‌ها در دو مرحله آموزش

جدول ۱. ارزیابی دقت ۴ نوع توابع هسته در طبقه‌بندی دمای متوسط ایستگاه سینوپتیک شیراز

مرحله آزمایش			مرحله آموزش			نوع الگو
MAE	RMSE	R ²	MAE	RMSE	R ²	
۰/۵۲	۰/۷۲	۰/۹۹	۰/۴۶	۰/۶۸	۰/۹۹	خطی
۰/۶۸	۰/۷۸	۰/۹۸	۰/۷۹	۰/۸۵	۰/۹۷	تابع پایه شعاعی
۰/۴۶	۰/۶۸	۰/۹۸	۰/۷۲	۰/۷۲	۰/۹۸	چندجمله‌ای
۰/۸۶	۰/۹۸	۰/۹۷	۰/۸۷	۰/۹۱	۰/۹۸	سیگموئید

جمع‌بندی و توصیه‌ها

برای بهبود عملکرد الگوهای ماشین‌های بردار پشتیبان در

پیش‌بینی اقلیمی، در ادامه چند توصیه ارائه شده است:

- ویژگی‌های یک ساختار ماشین‌های بردار پشتیبان به انتخاب نوع تابع هسته وابسته است. بنابراین، در انتخاب نوع تابع هسته دقت شود؛

- تعیین مؤلفه اصلی، فرآیندی مهم در تحلیل داده‌ها است و کمک می‌کند تا مؤلفه‌های اصلی تحقیق شناسایی و عوامل غیرضروری حذف شوند. بنابراین، باید در تعیین مؤلفه اصلی در پیش‌بینی اقلیمی تلاش شود تا با تعداد مؤلفه‌های بهینه بهترین برازش و پیش‌بینی به دست آید؛

- الگوهای مختلف ماشین‌های بردار پشتیبان مقایسه شده و عملکرد آن‌ها ارزیابی شود تا بهینه‌ترین ساختار برای پیش‌بینی مؤلفه‌های اقلیمی شناسایی شوند.

منابع

خاشعی سیوکی، عباس، شهیدی، علی، پور رضا بیلندی، محسن، امیرآبادی زاده، مهدی و جعفر زاده، احمد. (۱۳۹۷). بررسی عملکرد روش‌های ANN و SVR در ریزمقیاس نمایی بارش روزانه مناطق خشک. تحقیقات آب و خاک ایران (علوم کشاورزی ایران)، ۴۹(۴). صفوی، ح. (۱۳۸۸). هیدرولوژی مهندسی. انتشارات ارکان دانش. چاپ اول، ۶۰۷ ص.

پیش‌بینی‌ها کمک می‌کنند تا تغییرات آب و هوایی و الگوهای طبیعت بهتر درک شوند. در نتیجه، برای حفاظت از منابع طبیعی و جامعه بشری در برابر تغییرات آب و هوایی شدید، می‌توان راهبردهای مناسبی را برنامه‌ریزی کرد. پژوهش‌های مرتبط با پیش‌بینی تغییرات اقلیم، درک و عملکرد بهتری از تغییرات آینده را به دست می‌دهند. با استفاده از داده‌های جمع‌آوری‌شده و الگوریتم‌های پیچیده، پژوهشگران می‌توانند الگوهای تغییرات آب و هوایی را شبیه‌سازی کرده و نتایج را برآورد کنند. این نتایج به سازمان‌ها، سیاست‌گذاران و جامعه کمک می‌کند تا در حوزه‌های مختلف تصمیمات مناسبی بگیرند. در پیش‌بینی، طبقه‌بندی اطلاعات اهمیت زیادی دارد. با طبقه‌بندی اطلاعات می‌توان الگوها و روابط مختلف در داده‌ها را شناسایی و به تحلیل و پیش‌بینی بهتری دست یافت. در طول فرآیند طبقه‌بندی، اطلاعات مشابه گروه‌بندی شده و به‌عنوان چارچوب اصلی برای پیش‌بینی استفاده می‌شود. نتایج مربوط به الگوسازی تحت الگوهای ماشین‌های بردار پشتیبان نشان‌دهنده دقت مناسب این الگوها در پیش‌بینی مؤلفه‌های اقلیمی است.

- Chen, H. Wu, W. & Liu, H. B. (2016). Assessing the relative importance of climate variables to rice yield variation using support vector machines. *Theoretical and Applied Climatology*, 126, 105-111.
- Gocic, M. Shamshirband, S. Razak, Z. Petkovi, D. Ch, S. & Trajkovic, S. (2016). Long-term precipitation analysis and estimation of precipitation concentration index using three support vector machine methods. *Advances in Meteorology*, 2016.
- Guobin, F. Stephen, P. C. Francis, H.S.C. Jin, T. Hongxing, Z. Andrew, J. F. Wenbin, L. Sergey, K. 2013. Modelling runoff with statistically downscaled daily site, gridded and catchment rainfall series. *Journal of Hydrology* 492 (2013) 254-265.
- Haiyun Panda D K. Mishra A. Jena S K. James B K. Kumar A. 2007: *The Influence of Drought and Anthropogenic Effects on Groundwater Levels in Orissa, India*. 140-153.
- Hastie, Trevor, Robert Tibshirani, Jerome H. Friedman, *the Elements of Statistical Learning*, Springer. (2001). ISBN 0-387-95284-5.
- Kalra, A. Ahmad, S. (2012). Estimating Annual precipitation for the Colorado River Basin using oceanic-atmospheric oscillations. *Water Resources Research*, 48(6)
- Lu, Y. & Qin, X. S. (2014). A coupled K-nearest neighbour and Bayesian neural network model for daily rainfall downscaling. *International Journal of Climatology*, 34(11), 3221-3236.
- Min H. Jae, & Lee C. Young (2005). Bankruptcy prediction using support vector machine with optimal choice of kernel function parameters. *Expert Systems with Applications*, 28, 603-614.
- Mirarabi, A. Nassery, H. R. Nakhaei, M. Adamowski, J. Akbarzadeh, A. H. & Alijani, F. (2019). Evaluation of data-driven models (SVR and ANN) for groundwater-level prediction in confined and unconfined systems. *Environmental Earth Sciences*, 78(15), 1-15.
- Nsubuga, F.W.N. Botai, O.J. Olwoch, J.M. Dew Rauten bach, C.J.Yvette, B. & Adebayo, O.A. (2014). The nature of rainfall in the main drainage sub-basins of Uganda. *Hydrological Sciences Journal*, 59 (2), 278-299. DOI: 10.1080/02626667.2013.804188.
- Sachindra, D. A. Huang, F. Barton, A. & Perera, B. J. C. (2013). Least square support vector and multi-linear regression for statistically downscaling general circulation model outputs to catchment streamflows. *International Journal of Climatology*, 33(5), 1087-1106.
- Shin S. Kyung, Lee S. Taik, & Kim J. Hyun (2005). An application of support vector machines in bankruptcy prediction model. *Expert Systems with Applications*, 28, 127-135.
- Stuart J, Russell, Peter Norvig, *Artificial Intelligence: A Modern Approach (Third Ed)*. Prentice Hall. (2010). ISBN 9780136042594.
- Tripathi, S. Srinivas, V. V. & Nanjundiah, R. S. (2006). Downscaling of precipitation for climate change scenarios: a support vector machine approach. *Journal of hydrology*, 330(3-4), 621-640.
- Vapnik, V. (1995). *The nature of statistical learning theory*. NY: Springer-Verlag.
- Vapnik, Vladimir, 2000, *The nature of statistical learning theory*. Springer Science & Business Media.
- Wen, X. Si, J. He, Z. Wu, J. Shao, H. & Yu, H. (2015). Support-vector-machine-based models for modeling daily reference evapotranspiration with limited climatic data in extreme arid regions. *Water resources management*, 29, 3195-3209.

The use of support vector machines in classification of climatic data

Seyed Masoud Soleimanpour¹, Parsa Haghighi²

1. Associate Professor, Soil Conservation and Watershed Management Research Department, Fars Agricultural and Natural Resources Research and Education Center, Agricultural Research, Education and Extension Organization (AREEO), Shiraz, Iran. (Corresponding author). Email: m.soleimanpour@areeo.ac.ir

2. Masters, Soil Conservation and Watershed Management Research Department, Fars Agricultural and Natural Resources Research and Education Center, Agricultural Research, Education and Extension Organization (AREEO), Shiraz, Iran Email: p.haghighi@areeo.ac.ir

Abstract

Identifying, predicting and managing crisis in a climate structure is of great importance. Models are used as practical tools for understanding complex systems and simulating and predicting their behavior. Support vector machines are one of the supervised learning methods used for classification and regression. Support vector machines are able to detect hidden patterns and respond to complex changes in climate data. In this article, the structure of the support vector machine method and its application in climate data classification are presented. The characteristics of the structure of support vector machines are related to the selection of the kernel function type, so sufficient care must be taken in the selection of the kernel function type and on the other hand, PCI in climate forecasting is an important step in climate forecasting in order to make the best fit between forecasting and predicted data with the optimal number of parameters.

Keywords: Support vector machines, Classification of climatic data, kernel functions, PCI.