



فصلنامه راهبرد مدیریت مالی

دانشگاه الزهرا

سال دهم، شماره سی و هشتم، پاییز ۱۴۰۱

صفحات ۲۶-۱



### مقاله پژوهشی

## پیش بینی شاخص کل بورس اوراق بهادار تهران با استفاده از رگرسیون بردار پشتیبان بر مبنای تکنیک کاهش ابعاد<sup>۱</sup>

سمیه محبی<sup>۲</sup>، محمداسماعیل فدائی نژاد<sup>۳</sup>، محمد اصولیان<sup>۴</sup>

تاریخ پذیرش: ۱۴۰۱/۰۶/۱۹

تاریخ دریافت: ۱۴۰۰/۰۱/۱۰

### چکیده

بازارهای سهام نقش مهمی در سازماندهی سیستم‌های اقتصادی مدرن دارند. پژوهش‌های گسترده‌ای در زمینه پیش‌بینی آن‌ها با استفاده از تکنیک‌های هوشمند انجام شده است. با توجه به این‌که دقت عملکرد این تکنیک‌ها به میزان قابل توجهی تحت تأثیر ویژگی‌های ورودی آن است، یکی از پیشرفت‌های به‌کار رفته در استفاده از مدل‌های هوشمند، علاوه بر کاربرد مدل‌های ترکیبی، استفاده از کاهش ابعاد به‌عنوان یک پیش‌مرحله برای مدل پیش‌بینی می‌باشد. در این پژوهش برای پیش‌بینی روزانه شاخص کل بورس اوراق بهادار تهران همزمان از دو روش کاهش ابعاد (انتخاب و استخراج) به منظور انتخاب ویژگی‌های مناسب به عنوان ورودی‌های مدل استفاده می‌شود. به‌طوری‌که برای انتخاب ویژگی‌ها از الگوریتم *mRMR-MID* و برای استخراج ویژگی‌ها از الگوریتم *PCA* استفاده می‌شود. سپس از رگرسیون بردار پشتیبان به‌عنوان مدل پیش‌بینی استفاده می‌شود. با توجه به نتایج بدست آمده از تحلیل استفاده از تکنیک‌های کاهش ابعاد در مدل پیش‌بینی، در نهایت الگوریتمی برای انتخاب ویژگی‌های مناسب بر شاخص، تحت عنوان *ISF\_MID* پیشنهاد می‌شود. نتایج نشان می‌دهد که با روش پیشنهادی، می‌توان با ۷ ویژگی انتخابی به دقت بالایی در پیش‌بینی روزانه شاخص کل بورس اوراق بهادار تهران با درصد خطا ۳/۴۶ دست یافت. لازم به ذکر است مدل‌های مورد بررسی در مرحله پیاده‌سازی با روش اعتبارسنجی متقابل *k-fold* مورد ارزیابی قرار گرفتند. همچنین از معیارهای *MAE*، *MSE* و *RMSE* برای ارزیابی عملکرد مدل‌های مذکور استفاده می‌شود.

**واژگان کلیدی:** پیش‌بینی شاخص بورس، رگرسیون بردار پشتیبان، تکنیک کاهش ابعاد، انتخاب ویژگی، تجزیه و تحلیل مولفه‌های اصلی.

**طبقه‌بندی موضوعی:** *C38, C53, C61, G10, G17*

۱. کد DOI مقاله: 10.22051/JFM.2022.35543.2528

۲. دانشجوی دکتری، رشته مدیریت مالی، دانشکده مدیریت و حسابداری، دانشگاه شهید بهشتی تهران، ایران. E-mail: somayeh.mohebi@gmail.com

۳. دانشیار، گروه مدیریت مالی و بیمه، دانشکده مدیریت و حسابداری، دانشگاه شهید بهشتی، تهران، ایران. نویسنده مسئول. E-mail: m-fadaei@sbu.ac.ir

۴. استادیار، گروه مدیریت مالی و بیمه، دانشکده مدیریت و حسابداری، دانشگاه شهید بهشتی، تهران، ایران. E-mail: m\_osoolian@sbu.ac.ir

## مقدمه

بازارهای سهام به دلیل داشتن پتانسیل‌های بالای مالی، توجه بسیاری از سرمایه‌گذاران را به خود جلب نموده است (نواسالمی<sup>۱</sup>، ۲۰۲۰). از این رو پیش‌بینی این بازارها در میان تحلیل‌گران خبره و سرمایه‌گذاران از اهمیت بالایی برخوردار بوده و بسیار چالش برانگیز است (کومار، سارانجی و ورما<sup>۲</sup>، ۲۰۲۱). این مسئله عمدتاً به این دلیل است که بازار سهام یک سیستم بی‌نظم، پراشوب، ناپارامتریک، پویا، غیرخطی و نامانای می‌باشد (رافیوزامن<sup>۳</sup>، ۲۰۱۴).

برای پیش‌بینی و تحلیل بازارهای سهام از مدل‌های آماری و یادگیری ماشین استفاده شده است. مدل‌های آماری مانند میانگین متحرک خود رگرسیون انباشته<sup>۴</sup> (ARIMA) واریانس شرطی خودرگرسیون تعمیم یافته<sup>۵</sup> (GARCH)، و نوسانات تصادفی<sup>۶</sup> (SV)، هنوز هم برای پیش‌بینی سری‌های زمانی در اقتصاد استفاده می‌شوند، دلیل اصلی آن این است که خصوصیات ریاضی آن‌ها به خوبی درک شده است (الحق و همکاران<sup>۷</sup>، ۲۰۲۱). با این حال، روش‌های آماری بازارهای سهام را خطی، ثابت و نرمال فرض می‌کنند که عملکرد و کاربرد عملی آن‌ها را در این حوزه محدود می‌کند. مدل‌های یادگیری ماشین مانند ماشین بردار پشتیبان<sup>۸</sup> (SVM)، شبکه‌های عصبی مصنوعی<sup>۹</sup> (ANN)، سیستم‌های فازی<sup>۱۰</sup> (FIS) و الگوریتم ژنتیک<sup>۱۱</sup> (GA)، در سال‌های اخیر به منظور توضیح رفتار قیمت‌ها و پیش‌بینی آن‌ها در بازار سرمایه به دلیل تحلیل‌های چند متغیره بدون هیچ فرض از پیش تعیین شده‌ای با ویژگی غیرخطی و سهولت تعمیم، محبوبیت فراوانی یافته‌اند (ژونگ و انکه<sup>۱۲</sup>، ۲۰۱۷). از میان این روش‌ها، SVM که یکی از روش‌های یادگیری با نظارت است به طور گسترده‌ای برای پیش‌بینی قیمت سهام مورد استفاده قرار گرفته است. استفاده از SVM برای پیش‌بینی تغییر قیمت روزانه در شاخص قیمت سهام ترکیبی کره<sup>۱۳</sup> (KOSPI) که توسط کیم در سال ۲۰۰۳ انجام شد، هنوز هم یکی از پراستنادترین کارها در این زمینه است (بوستوس و پومارس-کویمایا<sup>۱۴</sup>، ۲۰۲۰). اخیراً هنریکه، سوبریرو و کیمورا<sup>۱۵</sup> (۲۰۱۹)، ۵۴۷ مقاله را بررسی کردند و دریافتند که در ۳۷٪ مطالعات بررسی شده از SVM برای پیش‌بینی قیمت سهام استفاده نموده‌اند.

1. Nevasalmi
2. Kumar, Sarangi & Verma
3. Rafiuzzaman
4. Autoregressive Integrated Moving Average
5. Generalized Autoregressive Conditional Heteroskedasticity
6. Stochastic Volatility
7. Ul Haq et al.
8. Support Vector Machine
9. Artificial Neural Network
10. Fuzzy Inference System
11. Genetic Algorithm
12. Zhong & Enke
13. Korean composite stock price
14. Bustos & Pomares-Quimbaya
15. Henrique, Sobreiro & Kimura

هرچند طبق مطالعات انجام شده، ماشین بردار پشتیبان عملکرد جالب توجه و مناسبی در پیش‌بینی و طبقه‌بندی دارد، اما دقت عملکرد آن به‌میزان قابل توجهی تحت تأثیر نوع و تعداد ویژگی‌های ورودی آن است. بنابراین کاهش تعداد ویژگی‌هایی که باید در آموزش ماشین بردار پشتیبان به‌کار گرفته شود، تأثیر به‌سزایی در افزایش دقت نتایج و کاهش هزینه دارد (لی<sup>۱</sup>، ۲۰۰۹، ژانگ و همکاران<sup>۲</sup> ۲۰۱۴). با توجه به این‌که بازارهای سهام تحت تأثیر عوامل مرتبط زیادی از جمله وضعیت اقتصادی، متغیرهای خاص صنعت، چشم‌انداز شرکت، نفوذ روانشناختی سرمایه‌گذاران و سیاست‌های دولت قرار دارند و بسیاری از آن‌ها به‌عنوان متغیرهای ورودی در طول توسعه یک سیستم پیش‌بینی بازار سهام ممکن است مورد استفاده قرار گیرند. بنابراین اگر از SVM انتظار می‌رود که پیش‌بینی دقیق و کارآمدی را انجام دهد، ضروری است تا ویژگی‌هایی که حاوی مفیدترین اطلاعات باشند به‌عنوان ورودی‌های SVM انتخاب شوند. این نوع گزینش وظیفه اصلی تکنیک کاهش ابعاد<sup>۳</sup> می‌باشد. کاهش ابعاد می‌تواند با دو شیوه متفاوت تحت عنوان انتخاب و استخراج ویژگی، اجرا شود (ژونگ و انکه، ۲۰۱۷)، این دو روش اهمیت قابل توجهی در یادگیری ماشین دارند، زیرا منجر به کاهش پیچیدگی‌های محاسباتی شده و به بهبود توانایی تعمیم الگوریتم طبقه‌بندی کمک می‌کند. علت این امر کاهش اندازه نمونه‌های آموزشی می‌باشد که خطر "بیش برآزش"<sup>۴</sup> را کاهش می‌دهد (کوالکانانت و همکاران<sup>۵</sup>، ۲۰۱۶).

با وجود این‌که شناسایی و انتخاب ویژگی‌های ورودی با دقت پیش‌بینی بالا، مدت‌ها است که یک موضوع پژوهش در داده‌کاوی می‌باشد (الحق و همکاران، ۲۰۲۱)، و محققین در آمار، علوم کامپیوتر و ریاضیات کاربردی در این زمینه سال‌های زیادی کار کردند و انواع مختلفی از تکنیک‌های کاهش ابعاد خطی و غیرخطی را توسعه دادند، با این حال مطالعات موجود اتکا به استفاده از یک نوع تکنیک کاهش ابعاد (انتخاب یا استخراج ویژگی) را نشان می‌دهند. این امر ممکن است برخی از مفروضات مهم در مورد عملکرد رگرسیون اصلی متصل به متغیرهای ورودی و خروجی را نادیده بگیرد. بنابراین در این مطالعه، همزمان از دو شیوه متفاوت کاهش ابعاد استفاده می‌شود. به‌طوری‌که در روش انتخاب ویژگی از معیار کمینه‌ی افزونگی بیشینه‌ی وابستگی<sup>۶</sup> mRMR استفاده می‌شود، که برخلاف پژوهش‌های صورت گرفته در این حوزه که از مقدار همبستگی ویژگی‌های ورودی نسبت به خروجی مدل برای انتخاب ویژگی‌های مؤثر استفاده می‌شود، این روش آماری؛ ویژگی‌های مؤثر را با توجه به بیشینه‌سازی معیار وابستگی آماری مجموعه ویژگی‌ها با ویژگی هدف، و کمینه کردن اطلاعات متقابل در بین مجموعه ویژگی‌های انتخابی، گزینش می‌کند. برای استخراج ویژگی‌ها از تجزیه و تحلیل مولفه‌های اصلی<sup>۷</sup> (PCA) که محبوب‌ترین تکنیک خطی برای کاهش ابعاد می‌باشد، استفاده می‌شود. سپس با ارزیابی عملکرد دو نوع تکنیک کاهش ابعاد، الگوریتمی به‌منظور

1. Lee
2. Zhang et al.
3. Dimensionality Reduction
4. Overfit
5. Cavalcante et al.
6. Minimum Redundancy Maximum Relevance
7. Principal Component Analysis

انتخاب ویژگی‌های مناسب با هدف فیلتر کردن متغیرهای ورودی بی‌ربط و اضافی برای کاهش پیچیدگی و بهبود دقت مدل پیش‌بینی پیشنهاد می‌گردد. در نهایت با مقایسه عملکرد روش‌های مختلف داده‌کاوی، روش مناسب برای مدل پیش‌بینی انتخاب می‌گردد و بر اساس این روش ویژگی‌های مهم به‌عنوان ورودی-های مدل پیش‌بینی روزانه شاخص کل بورس اوراق بهادار تهران شناسایی می‌گردند. بنابراین برای نخستین-بار این پژوهش یک فرآیند داده‌کاوی جامع را برای پیش‌بینی شاخص روزانه بورس بر اساس ۶۹ ویژگی اقتصادی و مالی انجام می‌دهد. سپس یک مدل دو مرحله‌ای با بیشترین دقت برای پیش‌بینی روزانه شاخص کل بورس اوراق بهادار تهران ارائه می‌دهد.

## مبانی نظری و مروری بر پیشینه پژوهش

### مبانی نظری

پیشرفت سریع و گسترده علم و فناوری و تلاش محققان برای اصلاح روش‌ها و مدل‌های پیش‌بینی، اخیراً موجب افزایش شایان توجه گرایش به استفاده از تکنیک‌های پیشرفته هوشمند، به‌منظور توضیح رفتار قیمت‌ها و پیش‌بینی آن‌ها در بازار سرمایه شده است. این مدل‌ها به‌طور عمده از توانایی احصای فرایندهای غیر خطی، ناماننا و نویزی برخوردارند (باجلان، فلاح‌پور و دانا، ۱۳۹۶). یکی از مهم‌ترین روش‌های یادگیری هوشمند که در سال‌های اخیر در بازارهای مالی بسیار مورد استفاده قرار گرفته و به نتایج مطلوبی نیز دست یافته است، ماشین بردار پشتیبان (SVM) می‌باشد.

SVM ها متداول‌ترین الگوریتم در میان الگوریتم‌های جداسازی خطی به دو دلیل می‌باشند. (۱) داده‌ها در این مدل‌ها می‌توانند بدون داشتن فرضیات قوی به‌کار گرفته شوند. (۲) در حالی که بیشتر مدل‌های شبکه عصبی مرسوم ریسک تجربی (خطای طبقه بندی نادرست) را به‌حداقل می‌رسانند، SVM ریسک‌های ساختاری (خطای تعمیم) را نیز به‌حداقل می‌رساند (کاولکاننت و همکاران، ۲۰۱۶). وی (۲۰۱۲)، شهرت و کارآمدی مدل SVM را به واسطه فرمول ویژه تابع هدف محدب با محدودیت ضرایب لاگرانژ دانست. به‌گفته یوان (۲۰۱۳)، راه‌حل‌های SVM ممکن است در سطح جهانی بهینه باشد، در حالی که شبکه‌های عصبی مرسوم معمولاً راه‌حل بهینه محلی تولید می‌کنند حتی ممکن است در برخی مسائل محلی نیز با شکست روبه‌رو شوند.

SVM در واقع یک طبقه‌بندی‌کننده دو وجهی است که سعی دارد در دو طبقه، ابرصفحه‌ای ایجاد نماید که فاصله هر طبقه تا ابرصفحه حداکثر باشد. داده‌های نقطه‌ای که به ابر صفحه نزدیک‌ترینند، برای اندازه‌گیری این فاصله به‌کار می‌روند. این داده‌های نقطه‌ای، بردارهای پشتیبان نام دارند (منصورفر، غیور و خالق پرست اطهری، ۱۳۹۴). در این روش، ساخت مدل شامل دو مرحله آموزش و آزمایش می‌باشد. در انتهای فاز آموزش، قابلیت تعمیم مدل آموزش داده شده با استفاده از داده‌های آزمایش ارزیابی می‌شود (وی، ۲۰۱۲). از این رو، طراحی SVM کاملاً به استخراج زیر مجموعه‌ای از داده‌های آموزش (بردارهای پشتیبان) که توانایی نمایش دادن ویژگی‌های پایدار و ثابت داده‌ها را دارند، وابسته است (فلاح‌پور، گل ارضی و فتوره‌چیان، ۱۳۹۲).

با توجه به توانایی تخمین غیرخطی SVM، از آن هم در طبقه‌بندی<sup>۱</sup> (SVC) و هم مسائل رگرسیون<sup>۲</sup> (SVR) استفاده می‌شود (گوا-کیانگ<sup>۳</sup>، ۲۰۱۱). با استناد به نتایج پژوهش‌های بسیاری، می‌توان بیان کرد که SVR به‌عنوان گزینه جایگزین شبکه‌های عصبی مصنوعی و نیز روش‌های آماری، در کارهای تشخیص الگو و پیش‌بینی سری‌های زمانی مالی به طور گسترده‌ای مورد استفاده قرار گرفته است (کاوالکانت و همکاران، ۲۰۱۶). اما هنوز هم محدودیت‌هایی در الگوهای یادگیری SVR عمدتاً در حوزه مالی، که در آن‌ها داده‌ها بسیار نویزی، بی‌ثبات و دارای ابعاد بالا هستند، وجود دارد (کارا، بویاجی اوغلو و بایکان<sup>۴</sup>، ۲۰۱۱). بنابراین آماده‌سازی داده‌ها اولین قدم به‌منظور استفاده موفقیت‌آمیز هر روش هوشمند از جمله SVR است. به‌طوری‌که پس از تعریف متغیرهای ورودی و خروجی برای مدل‌سازی سری‌های زمانی مالی و به دست آوردن داده‌های ورودی مورد استفاده در آموزش، استفاده از برخی روش‌های پیش‌پردازش روی این داده‌ها ممکن است بسیار مفید باشد، زیرا می‌تواند منجر به بهبود دقت مدل یادگیری هوشمند گردد. یکی از مکانیسم‌های پیش‌پردازش داده‌ها، کاهش ابعاد می‌باشد. تعداد متغیرهای لازم برای اندازه‌گیری هر مشاهده اندازه "ابعاد داده" نامیده می‌شود. زمانی که داده‌ها دارای ابعاد زیاد هستند، علی‌رغم فرصت‌هایی که به‌وجود می‌آورند، ممکن است منجر به اطلاعات نامربوط و / یا زاید شوند، که نه تنها هزینه محاسباتی را افزایش می‌دهد بلکه می‌تواند عملکرد مدل پیش‌بینی را تضعیف نماید. بنابراین کاهش ابعاد تحلیل‌گران را قادر می‌سازد تا ویژگی‌های مناسب از یک مجموعه داده انتخاب و سپس از آن‌ها برای آموزش مدل پیش‌بینی استفاده کنند، که این امر می‌تواند منجر به کاهش پیچیدگی در مدل یادگیری و در نتیجه باعث بهبود عملکرد تعمیم الگوریتم طبقه‌بندی و دقت پیش‌بینی شود (کاوالکانت و همکاران، ۲۰۱۶).

کاهش ابعاد می‌تواند با دو شیوه متفاوت اجرا شود: ۱) با انتخاب مرتبط‌ترین متغیرها از مجموعه داده‌های اصلی که انتخاب ویژگی<sup>۵</sup> نامیده می‌شود، ۲) با تولید یک گروه کوچکتر از متغیرهای جدید، که هرکدام ترکیب مشخصی از متغیرهای ورودی قدیمی‌تر می‌باشند، که استخراج ویژگی<sup>۶</sup> نامیده می‌شوند (ژونگ و انکه، ۲۰۱۷). به‌گفته وب (۲۰۰۳)، یک تمایز مهم در مورد انتخاب ویژگی و استخراج ویژگی وجود دارد. مکانیسم‌های انتخاب ویژگی که به آن‌ها انتخاب زیر مجموعه ویژگی نیز گفته می‌شود، متغیرهایی را که برای مدل‌سازی داده‌های آموخته شده نامربوط هستند، شناسایی می‌کنند، و برخلاف روش‌های مبتنی بر استخراج ویژگی، این نوع روش‌ها معنای اصلی ویژگی‌ها را بعد از کاهش حفظ می‌کنند. در صورتی‌که، مکانیسم‌های استخراج ویژگی سعی می‌کنند مجموعه ویژگی‌های اصلی را به فضای ویژگی‌هایی با ابعاد پایین انتقال دهند بدون این‌که ماهیت مسئله تغییر کند (کاوالکانت و همکاران، ۲۰۱۶).

1. Classification
2. Regression
3. Guo-Qiang
4. Kara, Boyacioglu & Baykan
5. Feature Selection
6. Feature Extraction



انتخاب ویژگی، یکی از تکنیک‌هایی است که در مبحث یادگیری ماشین و همچنین شناسایی الگوهای آماری مطرح است. انتخاب ویژگی با تعیین زیرمجموعه‌ای از ویژگی‌های موجود که بیشترین اهمیت را برای مساله طبقه‌بندی دارند، ابعاد مجموعه اولیه ویژگی‌ها را کاهش می‌دهد (لیو و ژنگ<sup>۱</sup>، ۲۰۰۶). برای انتخاب ویژگی، راه‌حل‌ها و الگوریتم‌های فراوانی ارائه شده است. تمامی روش‌ها تلاش می‌کنند مجموعه‌ای از ویژگی‌ها را انتخاب کنند که علاوه بر توصیف کارآمد داده‌های ورودی، در عین حفظ دقت نتایج پیش‌بینی، متغیرهای نویزی و نامربوط را کاهش دهند. روش‌های مختلف انتخاب ویژگی را می‌توان به سه دسته؛ جستجوی کامل، جستجوی تصادفی و جستجوی مکاشفه‌ای<sup>۲</sup> تقسیم‌بندی نمود. در جستجوی کامل تمام زیرمجموعه‌های ممکن برای یافتن جواب بهینه بررسی می‌شوند، بنابراین این احتمال وجود دارد که بر روی مجموعه داده‌های بزرگ غیرعملی باشد، زیرا بسیار زمان‌بر می‌باشد. در جستجوی تصادفی محدوده کمتری از فضای کل حالات بررسی می‌شود. در این روش‌ها پیدا شدن جواب بهینه به اندازه منابع موجود و زمان اجرای الگوریتم بستگی دارد. در جستجوی مکاشفه‌ای از بحث همبستگی بین ویژگی‌های ورودی با یکدیگر و نیز با ویژگی خروجی، برای انتخاب ویژگی‌ها استفاده می‌کنند، که در علم آمار با دو اصطلاح توضیح داده می‌شود. اولی بحث ارتباط<sup>۳</sup> می‌باشد، که وابستگی بین ویژگی انتخابی به‌عنوان ورودی با خروجی مدل را نشان می‌دهد، و دومی افزونگی<sup>۴</sup> می‌باشد، که وابستگی بین ویژگی انتخابی به‌عنوان ورودی با ویژگی‌هایی که تا کنون به‌عنوان ورودی مدل انتخاب شده است، را نشان می‌دهد. در این راستا از یک الگو دقیق استفاده می‌شود که به‌حداکثر شدن وابستگی بین ویژگی ورودی با ویژگی خروجی، و حداقل شدن وابستگی بین ویژگی که قرار است انتخاب شود با ویژگی‌هایی که تاکنون انتخاب شده‌اند، منجر می‌گردد.

در روش استخراج ویژگی، فضای یک مسئله به فضای دیگر نگاشت می‌یابد. هدف از نگاشت، انتقال یک مسئله از فضایی که شامل تعداد زیادی ویژگی می‌باشد به فضایی که تعداد ویژگی‌های کمتری داشته باشد، است. در این نگاشت ماهیت مسئله به‌هیچ‌عنوان تغییر نمی‌کند و صرفاً در یک فضایی با ابعاد کمتر همان داده‌ها را خواهیم داشت. تجزیه و تحلیل مولفه‌های اصلی (PCA) قدیمی‌ترین و شناخته‌ترین روش آماری برای استخراج ویژگی‌های مهم از مجموعه داده‌هایی با ابعاد گسترده<sup>۵</sup> است. این متدولوژی به کار پیرسون<sup>۶</sup> (۱۹۰۱)، برمی‌گردد و مبتنی بر ایده تعریف یک سیستم یا فضای هماهنگ جدید است که داده‌های خام می‌توانند با استفاده از داده‌های بسیار کمتری بیان شوند، بدون آن‌که اطلاعات قابل توجهی از دست برود.

### پیشینه پژوهش

روند مطالعات در پژوهش‌های داخلی و خارجی، برتری روش ماشین بردار پشتیبان را به روش‌های پیشین، از جمله روش شبکه‌های عصبی و نیز روش‌های سنتی نشان می‌دهد. هرچند طبق مطالعات انجام شده، ماشین بردار پشتیبان، عملکرد جالب توجه و مناسبی در پیش‌بینی و طبقه‌بندی دارد، اما دقت عملکرد آن به‌میزان قابل توجهی

1. Lui & Zheng
2. Heuristic
3. Relevance
4. Redundancy
5. High- dimensional
6. Pearson

تحت تأثیر نوع و تعداد ویژگی ورودی آن است. به طوری که کاهش تعداد ویژگی‌هایی که باید در آموزش ماشین بردار پشتیبان به کار گرفته شود، تأثیر به‌سزایی در افزایش دقت نتایج و کاهش هزینه دارد. بنابراین، به تدریج روش‌های کاهش ابعاد برای مقابله با این محدودیت و بهبود عملکرد مدل‌های پیش‌بینی معرفی شدند. در ادامه به ذکر چند نمونه از مطالعاتی پرداخته می‌شود که اخیراً در داخل و خارج از کشور انجام گرفته و بیشترین ارتباط را با موضوع پژوهش دارند. در حوزه مطالعات پیش‌بینی روند، باجلان، فلاح‌پور و دانا (۱۳۹۶)؛ مدلی برپایه‌ی ماشین بردار پشتیبان وزن‌دهی شده همراه با روش انتخاب ویژگی هیبرید که مرکب از یک بخش فیلتر کننده و یک بخش پوشش دهنده به‌منظور انتخاب زیرمجموعه‌ای بهینه از ویژگی‌ها می‌باشد، ارائه نمودند. براساس نتایج، این پژوهش نشان داد مدل VW-SVM همراه با انتخاب ویژگی F-SSFS عملکرد بهتری در پیش‌بینی قیمت سهم، نسبت به روش‌های موجود دارد.

فلاح پور، گل ارضی و فتوره چیان (۱۳۹۲) در پژوهشی تحت عنوان «پی‌بینی روند حرکتی قیمت سهام با استفاده از ماشین بردار پشتیبان برپایه الگوریتم ژنتیک در بورس اوراق بهادار تهران» اقدام به پیش‌بینی روند حرکت و تغییرات قیمت سهام در بازار بورس اوراق بهادار تهران نمودند. آن‌ها متغیرهای ورودی ماشین بردار پشتیبان را توسط الگوریتم ژنتیک، بهینه‌سازی نمودند. نتایج پژوهش آن‌ها نشان داد که ماشین بردار پشتیبان برپایه الگوریتم ژنتیک، دقت بسیار بیشتری در پیش‌بینی نسبت به ماشین بردار پشتیبان ساده دارد.

راعی، نیک عهد قصیرائی و حبیبی (۱۳۹۵)، در پژوهشی به‌منظور افزایش دقت پیش‌بینی شاخص بورس اوراق بهادار تهران ترکیبی از روش‌های آماری و هوش مصنوعی را به کار برده‌اند. در این پژوهش ابتدا از روش PCA برای پالایش اولیه داده‌ها استفاده شده است. سپس با استفاده از رگرسیون بردار پشتیبان بهینه شده به‌وسیله الگوریتم حرکت تجمعی ذرات، به پیش‌بینی شاخص اقدام شده است. نتایج بدست آمده نشان داد که پیش‌پردازش روی داده‌ها، خطای پیش‌بینی مدل را به‌طور قابل ملاحظه‌ای کاهش داده است. نگوین و لو<sup>۱</sup> (۲۰۱۴)، با ترکیب SOM<sup>۲</sup> که یک الگوریتم خوشه‌بندی است، و f-SVM<sup>۳</sup> به پیش‌بینی قیمت سهام پرداختند. با توجه به نتایج به‌دست آمده مدل دو مرحله‌ای یادشده نسبت به مدل‌های دیگر مانند مدل RBN، مدل ترکیبی SOM با SVM و مدل ANFIS عملکرد بهتری داشته است. لی پینگ نی، ژی وی نی و گائو<sup>۴</sup> (۲۰۱۱) در پژوهشی با عنوان "پیش‌بینی روند سهام بر اساس انتخاب ویژگی فراکتال<sup>۵</sup> و ماشین بردار پشتیبان؛ برای پیش‌بینی یک روز آتی شاخص قیمتی سهام، از رویکرد ترکیبی SVM با روش انتخاب ویژگی فراکتال که برای حل مشکلات غیرخطی مناسب است و در انتخاب بهینه تعداد ویژگی‌ها کمک می‌کند، استفاده کردند. نتایج نشان از عملکرد مناسب‌تر این روش نسبت به روش‌های مشابه داشته است.

1. Nguyen & Le
2. Self-Organizing Map (SOM)
3. Fuzzy -Support Vector Machines
4. Li-Ping Ni, Zhi-Wei Ni & Ya-Zhuo Gao
5. Fractal Feature Selection

او و وانگ<sup>۱</sup> (۲۰۰۹)، برای پیش‌بینی روند حرکتی قیمت شاخص بازار سهام هنگ‌کنگ، از ۱۰ روش داده‌کاوی استفاده کردند. این روش‌ها شامل تجزیه و تحلیل افتراقی خطی، تحلیل افتراقی درجه دوم، الگوریتم K-نزدیکترین همسایه، الگوریتم دسته‌بندی بیز براساس تابع برآورد کرنل، مدل لاجیت، طبقه‌بندی درختی، شبکه عصبی، طبقه‌بندی بیزی با فرایند گاوسی، ماشین بردار پشتیبان و حداقل مربعات ماشین بردار پشتیبان می‌باشند. نتایج تجربی نشان داد که عملکرد پیش‌بینی مدل‌های ماشین بردار پشتیبان و حداقل مربعات ماشین بردار پشتیبان بهتر از سایر مدل‌ها می‌باشد.

هوانگ<sup>۲</sup> (۲۰۱۲)، از GA برای بهینه‌سازی پارامترهای SVR و انتخاب ورودی‌های مدل استفاده کرد. ژانگ و همکاران (۲۰۱۴) از یک روش انتخاب ویژگی تحت عنوان CFS<sup>۳</sup> همراه با ۷ مدل پایه‌ای پیش‌بینی و ۱۸ ویژگی ورودی، برای پیش‌بینی روند تغییرات شاخص بورس شانگهای استفاده کردند. نتایج نشان داد که روش انتخاب ویژگی پیشنهادی آن‌ها از لحاظ دقت، عملکرد بسیار بهتری نسبت به سه روش انتخاب ویژگی پرکاربرد (PCA, CART, LASSO) دارد. سینگ و سریواستاوا<sup>۴</sup> (۲۰۱۷) از PCA برای سرعت بخشیدن به آموزش مدل بدون از دست دادن دقت در صحت پیش‌بینی استفاده کردند. در این مقاله، آن‌ها قیمت سهام را با استفاده ترکیب دو تکنیک قدرتمند PCA و شبکه عصبی عمیق پیش‌بینی کردند، سپس، نتایج خود را با شبکه عصبی تابع پایه شعاعی مقایسه کردند و نشان دادند دقت مدل پیشنهادی ۴۰.۸٪ بهبود یافته است.

ژانگ و انکه (۲۰۱۷) به منظور پیش‌بینی بازده شاخص بورس، سه تکنیک کاهش ابعاد شامل تحلیل مولفه اصلی (PCA)، تحلیل مؤلفه اصلی مقاوم فازی (FRPCA) و تحلیل مولفه اصلی مبتنی بر کرنل (KPCA) را بر ۶۰ داده مالی و اقتصادی به منظور ساده‌سازی و آرایش مجدد ساختار اصلی مورد استفاده قرار دادند. آن‌ها نتیجه گرفتند که استفاده از PCA نه تنها مقدار داده‌های مورد نیاز برای آموزش مدل‌ها را ساده می‌کند، بلکه باعث افزایش دقت کلی پیش‌بینی‌ها می‌شود.

### پرسش‌های پژوهش

این پژوهش به دنبال پاسخگویی به پرسش‌های زیر است:

۱. کدام یک از تکنیک‌های کاهش ابعاد (انتخاب یا استخراج ویژگی) تاثیر بیشتری در دقت عملکرد مدل SVR به منظور پیش‌بینی شاخص بورس اوراق بهادار تهران دارند؟
۲. برای تخمین mRMR کدام یک از الگوریتم‌های انتخاب ویژگی با نام‌های MID و FCD، مناسب می‌باشند؟
۳. الویت ویژگی‌های انتخاب شده توسط MID و FCD، به چه صورت است؟
۴. آیا الگوریتم پیشنهادی در این پژوهش در مقایسه با روش mRMR و PCA، عملکرد بهتری در انتخاب ویژگی‌های مناسب برای پیش‌بینی شاخص بورس اوراق بهادار تهران دارد؟

1. Ou & Wang
2. Huang
3. Causal Feature Selection
4. Singh & Srivastava



۵. چه ویژگی‌هایی برای ورودی مدل پیش‌بینی شاخص روزانه بورس اوراق بهادار تهران مناسب هستند و منجر به بهبود عملکرد مدل پیش‌بینی می‌شوند؟

### روش‌شناسی پژوهش

در این بخش از پژوهش به صورت اجمالی مجموعه داده، آماده‌سازی مجموعه‌ی داده، تکنیک‌های مختلف کاهش ابعاد، رگرسیون بردار پشتیبان به عنوان مدل پیش‌بینی شاخص بورس و معیارهای ارزیابی عملکرد مدل پیش‌بینی، معرفی می‌شوند.

### ایجاد مجموعه داده

نظر به اینکه افزایش تعداد داده‌ها در شبکه‌های عصبی مصنوعی موجب کسب نتایج دقیق‌تر می‌شود، لذا داده‌های مورد استفاده در پژوهش که از بررسی مطالعات پیشین و همچنین مبانی تئوریک استخراج گردیده‌اند، بصورت روزانه از تاریخ ۱۳۹۲/۱۰/۲۸ تا ۱۳۹۷/۵/۳۰ به مدت ۱۱۰۸ روز جمع‌آوری شدند. برای جمع‌آوری داده‌ها از پایگاه‌های اینترنتی شرکت خدمات فناوری بورس تهران، بانک مرکزی جمهوری اسلامی ایران، بانک داده‌های اقتصادی و مالی مربوط به وزارت امور اقتصادی و دارایی و اوپک بهره گرفته شده است.

**جدول ۱. لیست ویژگی‌های مورد بررسی برای پیش‌بینی شاخص کل بورس اوراق بهادار تهران**

نام ویژگی	توصیف مختصر ویژگی
Tehran Exchange Price Index (TEPIX)	بازده شاخص کل در روز جاری و سه روز گذشته
Trading Volume (TV)	تغییر نسبی حجم معاملات درسه روز گذشته
Index of 50 more active Companies (IC50)	بازده شاخص ۵۰ شرکت فعال بورس درسه روز گذشته
Price Index of 50 more active Companies (PIC50)	بازده شاخص قیمت ۵۰ شرکت فعال بورس درسه روز گذشته
30 Largest Companies Index (LCI30)	بازده شاخص ۳۰ شرکت بزرگ درسه روز گذشته
Industrial Index (InI)	بازده شاخص صنعت درسه روز گذشته
Financial Index (FI)	بازده شاخص مالی درسه روز گذشته
Exponential Moving Average (EMA)	میانگین متحرک نمایی قیمت پایانی ۱۰، ۲۰، ۵۰ و ۲۰۰ روز گذشته
Moving Average (MA)	میانگین متحرک قیمت پایانی ۵، ۱۰، ۳۰ روز گذشته
Relative Difference in Percentage (RDP)	تفاوت نسبی در درصد بازده شاخص در ۵، ۱۰، ۱۵ و ۲۰ روز گذشته
Moving Price level Percentage (MPP)	شاخص سطح مقاومت قیمت <sup>۱</sup> برای ۳۰ و ۱۲۰ روز گذشته
Price / Earning per share (P/E)	نسبت قیمت به سود با ۴ وقفه و تغییرات نسبی ۳ روز گذشته آن
OPEC Oil Prices (OP)	قیمت نفت اوپک با ۴ وقفه و تغییرات نسبی ۳ روز گذشته آن
USD/IRR (DP)	قیمت دلار با ۴ وقفه و تغییرات نسبی ۳ روز گذشته آن
Euro/IRR (EP)	قیمت یورو با ۴ وقفه و تغییرات نسبی ۳ روز گذشته آن
Gold Coin Price (GP)	قیمت سکه طلا با ۴ وقفه و تغییرات نسبی ۳ روز گذشته آن

۱. مفهومی در تحلیل تکنیکی است که گرایش سهام و شاخص‌ها را در پشتیبانی از کاهش قیمت یا مقاومت در برابر افزایش قیمت، ارزیابی می‌کند. MPP بالا بیان‌گر سطح بیشتر مقاومت می‌باشد.

## آماده‌سازی مجموعه داده

برای یکسان کردن مقیاس پارامترها، مرحله نرمال‌سازی انجام شده است. برای نرمال‌سازی پارامترها که شامل تعیین مقدار هر پارامتر در بازه‌ی بین صفر و یک است، از رابطه ساده ۱ استفاده می‌شود:

$$\eta_{i,j} = \frac{v_{i,j} - \min(V_j)}{\max(V_j) - \min(V_j)} \quad (1)$$

در این رابطه،  $V_j$  مجموعه‌ی مقادیر پارامتر  $j$ ام و  $v_{i,j} \in V_j$  که  $i$  شماره نمونه در مجموعه داده است.  $\eta_{i,j}$  مقدار نرمال شده نمونه  $i$ ام از پارامتر  $j$  است.

## اولویت‌بندی ویژگی‌های مؤثر بر شاخص

در یادگیری ماشین، تعداد زیاد ورودی‌های مدل، باعث کاهش تعمیم مدل و افزایش سربار محاسباتی آن می‌شود. کاهش تعمیم، به معنی کاهش انطباق مدل با داده‌هایی است که آموزش مدل با آن‌ها انجام نشده است. بنابراین، کاهش تعداد ویژگی‌های ورودی، می‌تواند باعث افزایش دقت مدل پیش‌بینی گردد. الگوریتم‌های متعددی شامل الگوریتم‌های هوشمند و الگوریتم‌های مبتنی بر الگوهای آماری برای انتخاب ویژگی پیشنهاد شده‌اند. اغلب الگوریتم‌های آماری علاوه بر اینکه سرعت بالاتری در اجرا دارند می‌توانند الویت هر ویژگی را بصورت عددی مشخص کنند در صورتی که الگوریتم‌های هوشمند اغلب یک مجموعه از ویژگی‌ها را به‌عنوان ویژگی‌های مناسب انتخاب می‌کنند و ترتیبی از الویت‌ها ارائه نمی‌دهند. استفاده از مقدار همبستگی ویژگی‌های ورودی نسبت به خروجی مدل، یکی از روش‌های آماری انتخاب ویژگی برای مدل پیش‌بینی است. اما بهتر است انتخاب ویژگی بر اساس وابستگی بین داده‌ها انجام شود. بنابراین، به منظور انتخاب ویژگی‌های مناسب برای مدل پیش‌بینی معیاری نیاز است علاوه بر توجه به همبستگی هر ویژگی با مقدار شاخص آتی، همبستگی بین ویژگی‌های ورودی با یکدیگر نیز بررسی شود. در این راستا، ویژگی‌ها با توجه به بیشینه‌سازی معیار وابستگی آماری مجموعه ویژگی‌ها با ویژگی هدف و کمینه کردن اطلاعات متقابل (MI) در بین مجموعه ویژگی‌های انتخابی، گزینش می‌شوند. MI بین دو ویژگی  $X$  و  $Y$  با رابطه ۲ محاسبه می‌شود.

$$I(x; y) = \iint p(x, y) \log \frac{p(x)}{p(x)p(y)} dx dy \quad (2)$$

در رابطه فوق،  $p(x)$ ،  $p(y)$  و  $p(x, y)$  به ترتیب، توابع چگالی احتمال متغیرهای  $x$ ،  $y$  و وقوع هم‌زمان آن‌ها می‌باشند.

از رابطه ۳ برای محاسبه زیرمجموعه‌ی ویژگی‌ها با بیش‌ترین مقدار وابستگی با ویژگی هدف، استفاده می‌شود.

$$\max V_I(S, h), \quad V_I = \frac{1}{|S|} \sum_{i \in S} I(i, h) \quad (3)$$

$S$  زیرمجموعه‌ای از ویژگی‌های اولیه و  $V_I$  مقدار وابستگی  $S$  را با ویژگی هدف  $h$  نشان می‌دهد. انتخاب زیرمجموعه براساس بیش‌ترین وابستگی، می‌تواند شامل ویژگی‌هایی باشد که خود آن‌ها همبستگی بالایی باهم داشته باشند. بر این اساس، شرط حداقل افزونگی، برای یافتن زیرمجموعه‌ای از ویژگی‌ها دارای کمترین همبستگی با یکدیگر، ارائه گردیده است و از طریق رابطه ۴ محاسبه می‌شود.

$$\min W_I(S), \quad W_I = \frac{1}{|S|^2} \sum_{i,j \in S} I(i,j) \quad (4)$$

$W_I$  میانگین MI بین ویژگی‌ها در زیرمجموعه  $S$  است. با ترکیب دو رابطه ۳ و ۴، مفهومی تحت عنوان کمینه افزونگی بیشینه وابستگی (mRMR) معرفی گردیده است (پرز-رودریگز و همکاران، ۲۰۰۴). اساس mRMR با رابطه ۵ تعریف می‌شود:

$$\max \varphi(V, W), \quad \varphi = V - W \quad (5)$$

که  $\varphi(V, W)$  عملگری برای ترکیب روابط بیشینه‌ی وابستگی و کمینه‌ی افزونگی است. انتخاب فهرست بهترین ویژگی‌ها برای مجموعه داده‌ای با تعداد زیادی ویژگی، به دلیل نمایی بودن تعداد محاسبه‌های لازم برای یافتن حداکثر مقدار  $\varphi(V, W)$ ، غیرعملی است. به همین دلیل روش‌هایی برای تخمین mRMR معرفی گردیده‌اند. از متداول‌ترین آن‌ها، می‌توان به روش تفاضل اطلاعات متقابل<sup>۲</sup> (MID) اشاره کرد (دینگ و پنگ، ۲۰۰۵). MID با رابطه ۶ محاسبه می‌شود. باید توجه داشت که برای استفاده از MID نیاز به گسسته‌سازی داده‌ها است. بدین منظور با توجه به مقادیر بین کمینه و بیشینه‌ی شاخص بورس، تعداد بازه‌ها تعیین می‌شود.

$$\max_{i \in \varphi(S)} [I(i, h) - \frac{1}{|S|} \sum_{j \in S} I(i, j)] \quad (6)$$

همچنین، رابطه‌ی ۷ تحت عنوان FCD<sup>۴</sup> برای تخمین mRMR برای متغیرهای پیوسته ارائه شده است (ماندال و موخوپادیای، ۲۰۱۳،<sup>۵</sup>)

$$\max_{i \in \varphi(S)} [F(i, h) - \frac{1}{|S|} \sum_{j \in S} c(i, j)] \quad (7)$$

در رابطه فوق،  $c(i, j)$  مقدار همبستگی دو ویژگی  $i$  و  $j$  را نشان می‌دهد. تابع  $F$  با رابطه ۸ تعریف شده‌است.

$$F(g_i, h) = [\sum_k n_k (\bar{g}_k - \bar{g}) / (K - 1)] / \sigma^2 \quad (8)$$

مقدار تابع  $F$  برای متغیر  $g_i$ ، با توجه به  $K$  دسته‌ای که بر اساس  $h$  تعیین شده، محاسبه می‌شود.  $\bar{g}$  میانگین مقدار  $g_i$  در همه نمونه‌ها و  $\bar{g}_k$  میانگین مقدار در دسته  $k$ ام است. همچنین  $\sigma^2$  که واریانس تلفیقی است با رابطه ۹ محاسبه می‌شود.

$$\sigma^2 = [\sum_k (n_k - 1) \sigma_k^2] / (n - K) \quad (9)$$

در رابطه فوق،  $\sigma_k^2$  و  $n_k$  واریانس و اندازه دسته  $k$ ام است.

- 
1. Perez-Rodriguez et al.
  2. Mutual Information Difference
  3. Ding & Peng
  4. F-Test Correlation Difference
  5. Mandal & Mukhopadhyay

## تجزیه و تحلیل مولفه‌های اصلی (PCA)

تحلیل مولفه‌های اصلی، محبوب‌ترین تکنیک خطی برای کاهش ابعاد است، و به‌عنوان یکی از اولین تکنیک‌های چند متغیری، به‌دنبال ساخت یک نماینده با ابعاد کمتر از داده‌ها در عین حفظ واریانس حداکثری و ساختار کوواریانس داده‌ها است.

اساساً مولفه‌های اصلی، ترکیب‌های خطی از همه عوامل با ضرایب برابر با عناصر بردارهای ویژه<sup>۱</sup> می‌باشند. مقادیر مختلف مولفه‌های اصلی می‌توانند نسبت‌های مختلف ساختار واریانس-کوواریانس داده‌ها را توضیح دهند. مقادیر ویژه<sup>۲</sup> می‌توانند برای رتبه‌بندی بردارهای ویژه بر این اساس که چه مقدار از تغییرات داده‌ها توسط هر مولفه اصلی گرفته می‌شود، مورد استفاده واقع شوند. به‌طور دقیق، چنانچه  $\lambda^* = \{\lambda_i^*\}_{i=1}^M$  نشان‌دهنده مقادیر ویژه ماتریس همبستگی  $\text{COIT}(X)$  باشد به‌طوری که  $\lambda_1^* \geq \lambda_2^* \geq \dots \geq \lambda_M^*$  همچنین بردارهای  $e_i^T = (e_{i1}, e_{i2}, \dots, e_{iM})$  نشان‌دهنده بردارهای ویژه‌ای از  $\text{COIT}(X)$  مطابق با مقادیر ویژه  $\lambda_i^*, i = 1, 2, \dots, M$  می‌باشند. عناصر این بردارهای ویژه، ضرایب مولفه‌های اصلی هستند. یعنی، مولفه‌های اصلی داده‌های استاندارد شده عبارت است از  $Z = (Z_1, Z_2, \dots, Z_M)$  به طوری که؛

$$Z_W^T = (Z_{1W}, Z_{2W}, \dots, Z_{NW}), Z_{VW} = \frac{X_{VW} - \mu_W}{\sigma_W}, V = 1, 2, \dots, N \text{ and } W = 1, 2, \dots, M \quad (10)$$

بنابراین، رابطه ۱۰ می‌تواند به این شکل ذیل نوشته شود:

$$Y_i = \sum_{j=1}^M e_{ij} Z_j, \quad i = 1, 2, \dots, M \quad (11)$$

و اثبات شده است که:

$$\text{var}(Y_i) = \sum_{k=1, l=1}^M e_{ik} \text{corr}(X_k, X_l) e_{il} = e_i^T \rho e_i = \lambda_i^* \quad (12)$$

و

$$\text{cov}(Y_i, Y_j) = \sum_{k=1, l=1}^M e_{ik} \text{corr}(X_k, X_l) e_{jl} = e_i^T \rho e_j = 0 \quad (13)$$

در نهایت، با استفاده از برهان تجزیه طیفی، می‌توان گفت؛

$$\rho = \sum_{i=1}^M \lambda_i^* e_i e_i^T \quad (14)$$

و این واقعیت که هر دو  $e_i^T e_i = \sum_{j=1}^M e_{ij}^2 = 1$  و بردارهای ویژه مختلف نسبت به یکدیگر عمود هستند به‌طوری که  $e_i^T e_j = 0$ ، بنابراین واریانس  $i^{\text{th}}$  (بزرگترین) مولفه اصلی برابر با  $i^{\text{th}}$  بزرگترین مقدار ویژه است و مولفه‌های اصلی با یکدیگر ناهمبسته هستند. از آنجا که تغییرات کل  $Z$  به‌عنوان ردیابی ماتریس همبستگی  $\rho$  تعریف می‌شود، یعنی  $\text{trace}(\rho) = \sum_{i=1}^M \lambda_i^*$ ، نسب تغییرات توضیح شده توسط  $i^{\text{th}}$  مولفه اصلی به شکل  $\lambda_i^* / \text{trace}(\rho)$  تعریف می‌گردد که  $i = 1, 2, \dots, m$ . نسبت تغییرات توضیح شده توسط  $k$

1. Eigenvectors  
2. Eigenvalues

اول مولفه اصلی به عنوان مجموع  $k$  اول مقادیر ویژه که توسط  $trace(\rho)$  تفکیک شده، تعریف می شود.، یعنی  $\sum_{i=1}^k \lambda_i^* / \sum_{i=1}^M \lambda_i^*$ . به لحاظ نظری، اگر نسبت تغییراتی که توسط  $k$  اول مولفه های اصلی توضیح داده می شود بزرگ باشد، با کاهش ابعاد فضای داده از  $M$  به  $k$ ، اطلاعات زیادی از دست نمی رود. برای تعیین این که چه تعداد و کدام مولفه های اصلی باید به عنوان ورودی ها برای طبقه بندی مورد استفاده قرار گیرند، لازم است توازن بین دقت پیش بینی مورد انتظار، هزینه (زمان و ...) و پیچیدگی سیستم در نظر گرفته شود. به این معنی که مولفه های اصلی انتخابی، باید داده ها را به بهترین شکل، در حالی که ساختار داده ها را تا حد امکان ساده کرده است، توضیح دهند.

### رگرسیون بردار پشتیبان (SVR)

یکی از مدل های متداول برای توابع تخمین، رگرسیون بردار پشتیبان است، که برای حل مسائل رگرسیون غیرخطی مناسب است. در ادامه توضیحات مختصری در رابطه با نحوه عملکرد مدل ارائه می گردد. مجموعه داده  $M$  را در نظر بگیرید که شامل بردارهای ورودی  $x_i$  و خروجی متناظر  $y_i$  است. تعداد نمونه های این مجموعه داده برابر  $n$  است.

$$M = \left\{ (x_i, y_i) \mid i = 1, 2, \dots, n \right\} \quad (15)$$

در رگرسیون، به دنبال تخمین تابع  $f(x_i)$  هستیم که خروجی های آن حداقل فاصله را با مقادیر  $y_i$  داشته باشد. در رابطه ۱۶،  $\delta$  خطای تصادفی با توزیع  $N(0, \sigma^2)$  است.

$$y_i = f(x_i) + \delta \quad (16)$$

برای حل یک مسئله رگرسیون غیرخطی با SVR، ابتدا با بهره گیری از رابطه ۱۷، ورودی ها به صورت غیرخطی به فضای ویژگی  $f$  با ابعاد زیاد که به صورت خطی با خروجی وابستگی دارند، نگاشت می شوند. در این رابطه،  $w$  بردار وزن،  $b$  مقدار بایاس و  $\varphi(x)$  تابعی است که ورودی ها را از فضای  $R$  به فضای  $R^{N \times h}$  تصویر می کند.

$$f(x_i) = w\varphi(x_i) + b \quad |w \in R^{N \times h}, b \in R \quad (17)$$

تابع خطی  $f$  به گونه ای است که به اندازه  $\varepsilon$  از مقادیر واقعی انحراف دارد. برای داده هایی که خارج از باند باشند، از تابع جریمه ای استفاده می شود که با رابطه ۱۸ ارائه گردیده است.

$$L_\varepsilon(y_i, f(x_i)) = \begin{cases} 0 & |y_i - f(x_i)| \leq \varepsilon \\ |y_i - f(x_i)| - \varepsilon & \text{otherwise} \end{cases} \quad (18)$$

برای محاسبه ریسک عملیاتی تابع  $f$  از رابطه ۱۹ استفاده می شود.

$$R_{emp}[f] = \sum_1^n L_\varepsilon(y_i, f(x_i)) \quad (19)$$

برای کاهش ریسک عملیاتی مربوط به داده ها با توجه به بهینه سازی تابع رگرسیون از رابطه ۲۰ بهره گرفته می شود.

$$J = \frac{1}{2} \|w\|^2 + CR_{emp}[f] \quad (20)$$

برای داده‌هایی که مقدار  $|y - f(x_i)|$  آن‌ها بیشتر از  $\varepsilon$  باشد، مقدار  $\xi_i^+$  و یا  $\xi_i^-$  که مقدار تخطی را نشان می‌دهد، با رابطه‌های ۲۱ و ۲۲ محاسبه می‌شود.

$$\xi_i^+ = y - f(x_i) - \varepsilon \quad (21)$$

$$\xi_i^- = \varepsilon - y - f(x_i) \quad (22)$$

همچنین، رابطه‌ی ۲۳ بین تابع جریمه و مقدارهای تخطی وجود دارد.

$$L_\varepsilon(y_i, f(x_i)) = \xi_i^+ + \xi_i^- \quad (23)$$

برای تخمین تابع  $f$  از تابع هدف با رابطه‌ی ۲۴ استفاده می‌شود.

$$\text{minimize } \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n (\xi_i^+ + \xi_i^-) \quad (24)$$

S.t.  $\forall i$

$$-y_i + f(x_i) + \varepsilon + \xi_i^+ \geq 0$$

$$y_i - f(x_i) + \varepsilon + \xi_i^- \geq 0$$

$$\xi_i^+, \xi_i^- \geq 0$$

سپس، از ضریب لاگرانژ برای ایجاد شکل دوگان رابطه ۲۴ استفاده شده و ساده‌سازی انجام می‌شود. در صورتی که  $\alpha_i^+$  و  $\alpha_i^-$  به ترتیب ضریب قیدهای اول و دوم رابطه ۲۴ باشند، بعد از ساده‌سازی، رابطه ۲۵ حاصل می‌شود.

$$\text{minimize } \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n (\alpha_i^+ - \alpha_i^-) (\alpha_j^+ - \alpha_j^-) < \varphi(x_i) \cdot \varphi(x_j) > - \sum_{i=1}^n (\alpha_i^+ - \alpha_i^-) y_i + \varepsilon \sum_{i=1}^n (\alpha_i^+ + \alpha_i^-) \quad (25)$$

$$\text{S. t. } \left\{ \begin{array}{l} \sum_{i=1}^n (\alpha_i^+ - \alpha_i^-) = 0 \\ \alpha_i^+, \alpha_i^- \in [0, C] \end{array} \right\}$$

در مسائل غیرخطی، ضرب داخلی دو تابع  $\varphi(x_i)$  و  $\varphi(x_j)$  با تابع کرنل گوسی، در رابطه ۲۶ جایگزین می‌شود.

$$K(x_i, x_j) = \exp\left(-\frac{\|x_i - x_j\|^2}{2\sigma^2}\right) \quad (26)$$

در نهایت، تابع  $f$  با رابطه ۲۷ محاسبه می‌شود.

$$f(x) = \sum_{i=1}^n (\alpha_i^+ - \alpha_i^-) K(x_i, x) + b \quad (27)$$

در این رابطه، محاسبه  $b$  با فرمول ۲۸ انجام می‌شود.  $SV$  بردار پشتیبان می‌باشد.

$$b = \frac{1}{n} \left\{ \sum_{0 < \alpha_i^- < C} [y_i - \sum_{x_j \in SV} (\alpha_j^+ - \alpha_j^-) K(x_i, x_j) - \varepsilon] + \sum_{0 < \alpha_i^+ < C} [y_i - \sum_{x_j \in SV} (\alpha_j^+ - \alpha_j^-) K(x_i, x_j) + \varepsilon] \right\} \quad (28)$$



### معیارهای ارزیابی مدل‌های پیش‌بینی

روش‌های متنوعی برای ارزیابی مدل در مرحله‌ی پیاده‌سازی وجود دارد. این مقاله، از روش اعتبارسنجی متقابل k-fold<sup>۱</sup>، برای انتخاب مجموعه‌های آموزش و آزمایش بهره برده است. در این رویکرد k مرحله‌ای، مجموعه داده به ۱۰ دسته متمایز تقسیم‌بندی می‌شوند. در هر مرحله، ۱-۱۰ دسته برای آموزش مدل و یک دسته برای آزمایش مدل انتخاب می‌شوند. البته، پیش از تقسیم‌بندی داده‌ها، ابتدا جایگشتی بر روی مجموعه داده انجام شده تا نمونه‌ها به هم ریخته شوند. این به هم ریختگی باعث تنوع نمونه‌ها در توزیع داده‌ها می‌شود. در این پژوهش، برای بررسی دقت مدل‌های پیش‌بینی، از معیارهای میانگین قدر مطلق خطا<sup>۲</sup> (MAE)، میانگین مربع خطا<sup>۳</sup> (MSE)، و ریشه میانگین مربع خطا<sup>۴</sup> (RMSE) استفاده شده است. معیار MAE با رابطه‌ی ۲۹ محاسبه می‌شود.

$$MAE = \frac{1}{N} \sum_{i=1}^N |t_i - y_i| \quad (29)$$

N تعداد نمونه‌ها،  $t_i$  مقدار هدف و  $y_i$  مقدار پیش‌بینی شده برای نمونه‌ی  $i$ ام است. همچنین، MSE با رابطه‌ی ۳۰ محاسبه می‌شود.

$$MSE = \frac{1}{N} \sum_{i=1}^N (t_i - y_i)^2 \quad (30)$$

در نهایت، RMSE با رابطه‌ی ۳۱ محاسبه می‌گردد.

$$RMSE = \sqrt{\frac{1}{N} \sum_{i=1}^N (t_i - y_i)^2} \quad (31)$$

### تجزیه و تحلیل داده‌ها

#### انتخاب ویژگی‌های مؤثر بر شاخص با استفاده از روابط تخمین mRMR

از مزیت‌های استفاده از روابط تخمین mRMR می‌توان به این نکته اشاره کرد که با آن که تخمین قابل قبولی از mRMR دارند، ولی پیچیدگی محاسباتی آن‌ها کم است. به بیانی دیگر، سربار پایین، سرعت و قابلیت اطمینان بالا، باعث استفاده از این روش‌های تخمین برای انتخاب ویژگی‌های مناسب گردیده است. روش‌های آماری مورد استفاده برای تخمین mRMR در این مقاله، MID و FCD می‌باشند. بنابراین ویژگی‌های موجود در مجموعه داده ورودی را یک‌بار با استفاده از الگوریتم MID و بار دیگر با استفاده از FCD، در بردار  $\zeta$  بر اساس میزان وابستگی ویژگی‌ها با مقدار خروجی اولویت‌بندی می‌شوند. سپس در  $|\zeta|$  گام، زیرمجموعه‌هایی شامل  $|\zeta| \leq z \leq 1$  عنصر ابتدایی از  $\zeta$  را در مدل پیشنهادی ارزیابی کرده و میزان

1. K-fold Cross Validation
2. Mean Absolute Error
3. Mean Square Error
4. Root Mean Square Error



خطای پیش‌بینی برای هر یک اندازه‌گیری می‌گردد. ترتیب اولویت ویژگی‌ها با بهره‌گیری از MID و FCD به ترتیب در جدول‌های ۲ و ۳ نمایش داده شده است.

**جدول ۲. اولویت‌بندی ویژگی‌ها با الگوریتم MID با توجه به شاخص بورس یک روز آینده**

الویت	نام ویژگی	الویت	نام ویژگی	الویت	نام ویژگی	الویت	نام ویژگی	الویت	نام ویژگی
۱	RDP5	۱۵	IC50 <sub>2</sub>	۲۹	U <sub>3</sub>	۴۳	InI <sub>2</sub>	۵۷	GP <sub>d4</sub>
۲	TV <sub>1</sub>	۱۶	OP <sub>1</sub>	۳۰	OP <sub>d4</sub>	۴۴	MPP30	۵۸	P/E <sub>d2</sub>
۳	P/E <sub>d3</sub>	۱۷	PIC50 <sub>1</sub>	۳۱	P/E <sub>2</sub>	۴۵	D <sub>d3</sub>	۵۹	U <sub>d2</sub>
۴	D <sub>1</sub>	۱۸	U <sub>1</sub>	۳۲	PIC50 <sub>3</sub>	۴۶	TEPIX <sub>3</sub>	۶۰	D <sub>d1</sub>
۵	LCI30 <sub>1</sub>	۱۹	FI <sub>3</sub>	۳۳	U <sub>2</sub>	۴۷	EMA <sub>20</sub>	۶۱	OP <sub>d1</sub>
۶	RDP10	۲۰	OP <sub>3</sub>	۳۴	MPP120	۴۸	TEPIX <sub>2</sub>	۶۲	GP <sub>d3</sub>
۷	GP <sub>2</sub>	۲۱	P/E <sub>1</sub>	۳۵	InI <sub>1</sub>	۴۹	OP <sub>d3</sub>	۶۳	D <sub>d4</sub>
۸	OP <sub>2</sub>	۲۲	D <sub>3</sub>	۳۶	InI <sub>3</sub>	۵۰	P/E <sub>d4</sub>	۶۴	EMA10
۹	LCI30 <sub>3</sub>	۲۳	RDP20	۳۷	LCI30 <sub>2</sub>	۵۱	GP <sub>d1</sub>	۶۵	MA30
۱۰	GP <sub>3</sub>	۲۴	FI <sub>2</sub>	۳۸	GP <sub>d2</sub>	۵۲	U <sub>d1</sub>	۶۶	MA5
۱۱	FI <sub>1</sub>	۲۵	P/E <sub>3</sub>	۳۹	PIC50 <sub>2</sub>	۵۳	D <sub>d2</sub>	۶۷	MA10
۱۲	GP <sub>1</sub>	۲۶	TV <sub>3</sub>	۴۰	IC50 <sub>1</sub>	۵۴	P/E <sub>d1</sub>	۶۸	EMA50
۱۳	RDP15	۲۷	D <sub>2</sub>	۴۱	U <sub>d4</sub>	۵۵	U <sub>d3</sub>	۶۹	EMA200
۱۴	TV <sub>2</sub>	۲۸	TEPIX <sub>1</sub>	۴۲	IC50 <sub>3</sub>	۵۶	OP <sub>d2</sub>		

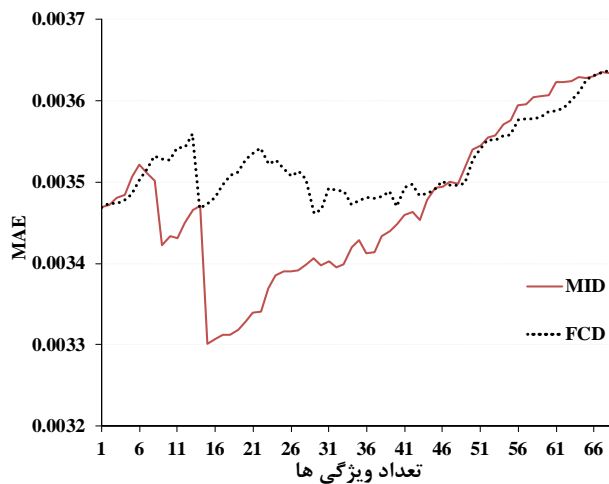
**جدول ۳. اولویت‌بندی ویژگی‌ها با الگوریتم FCD با توجه به شاخص بورس یک روز آینده**

الویت	نام ویژگی	الویت	نام ویژگی	الویت	نام ویژگی	الویت	نام ویژگی	الویت	نام ویژگی
۱	RDP5	۱۵	GP <sub>d2</sub>	۲۹	TEPIX <sub>3</sub>	۴۳	MA30	۵۷	D <sub>1</sub>
۲	RDP10	۱۶	GP <sub>d3</sub>	۳۰	LCI30 <sub>2</sub>	۴۴	FI <sub>3</sub>	۵۸	P/E <sub>d4</sub>
۳	IC50 <sub>1</sub>	۱۷	D <sub>d3</sub>	۳۱	IC50 <sub>2</sub>	۴۵	P/E <sub>2</sub>	۵۹	U <sub>1</sub>
۴	InI <sub>1</sub>	۱۸	D <sub>d4</sub>	۳۲	InI <sub>2</sub>	۴۶	GP <sub>3</sub>	۶۰	P/E <sub>d3</sub>
۵	PIC50 <sub>1</sub>	۱۹	U <sub>d1</sub>	۳۳	TEPIX <sub>2</sub>	۴۷	GP <sub>2</sub>	۶۱	P/E <sub>d1</sub>
۶	LCI30 <sub>1</sub>	۲۰	GP <sub>d4</sub>	۳۴	EMA200	۴۸	FI <sub>2</sub>	۶۲	GP <sub>1</sub>
۷	RDP15	۲۱	U <sub>d2</sub>	۳۵	U <sub>2</sub>	۴۹	D <sub>2</sub>	۶۳	OP <sub>3</sub>
۸	RDP20	۲۲	U <sub>d3</sub>	۳۶	P/E <sub>3</sub>	۵۰	OP <sub>1</sub>	۶۴	OP <sub>2</sub>
۹	P/E <sub>1</sub>	۲۳	PIC50 <sub>3</sub>	۳۷	MPP120	۵۱	OP <sub>d2</sub>	۶۵	TV <sub>2</sub>
۱۰	FI <sub>1</sub>	۲۴	U <sub>d4</sub>	۳۸	MA5	۵۲	OP <sub>d3</sub>	۶۶	P/E <sub>d1</sub>
۱۱	D <sub>d1</sub>	۲۵	PIC50 <sub>2</sub>	۳۹	EMA10	۵۳	OP <sub>d4</sub>	۶۷	TV <sub>3</sub>
۱۲	MPP30	۲۶	LCI30 <sub>3</sub>	۴۰	MA10	۵۴	U <sub>3</sub>	۶۸	TV <sub>1</sub>
۱۳	D <sub>d2</sub>	۲۷	IC50 <sub>3</sub>	۴۱	EMA20	۵۵	OP <sub>d1</sub>	۶۹	TEPIX <sub>1</sub>
۱۴	GP <sub>d1</sub>	۲۸	InI <sub>3</sub>	۴۲	EMA50	۵۶	D <sub>3</sub>		

لازم به ذکر است که به منظور تقسیم داده‌ها در دو دسته آموزش و آزمایش از روش اعتبارسنجی متقابل k-fold استفاده گردید. در آزمایش انجام‌شده مقدار K برابر با ۱۰ در نظر گرفته شده است. علت



انتخاب مقدار ۱۰ با توجه به متداول بودن استفاده از این عدد توسط محققان می‌باشد. در این حالت، هر نمونه از داده‌ها ۹ بار برای آموزش و یک‌بار برای آزمایش مدل مورد استفاده قرار می‌گیرد. مقدار MAE در نمودار ۱ برای پیش‌بینی شاخص بورس یک روز آینده، برای هر زیرمجموعه از ویژگی‌های انتخاب‌شده نمایش داده شده است.



نمودار ۱. استفاده از مدل SVR با ویژگی‌های اولویت‌بندی شده با روش‌های MID و FCD

با استفاده از این نمودار می‌توان مقایسه‌ای بین بهره‌گیری از MID و FCD در اولویت‌بندی ویژگی‌ها و اثر افزایش تعداد ویژگی‌های مدل پیش‌بینی داشت. یکی از نتایج قابل مشاهده در شکل این است که با افزایش تعداد ویژگی‌های اولویت‌بندی شده با FCD خطای مدل افزایش می‌یابد. در صورتی که با بهره‌گیری از اولویت‌بندی با MID، افزایش تعداد ویژگی‌ها ابتدا باعث افزایش دقت مدل می‌گردد ولی از ویژگی ۱۵ به بعد با کاهش دقت مدل روبه‌رو هستیم.

جدول ۴. مقایسه انتخاب ویژگی‌ها با استفاده از روش MID و FCD

FCD	MID	۱۵ ویژگی اول
۳/۷۳	۳/۵۵	درصد خطا
۰/۰۰۳۴	۰/۰۰۳۳	MAE
۰/۰۰۰۰۳۲	۰/۰۰۰۰۲۹	MSE
۰/۰۰۰۵۶	۰/۰۰۰۵۳	RMSE

نام ۱۵ ویژگی ابتدایی که توسط MID اولویت‌بندی شده و بالاترین دقت را در مدل SVR موجب می‌شوند، در جدول ۵ آمده است. از دیگر نکات قابل ملاحظه این است که اثر دو ویژگی ۹ و ۱۵ یعنی بازده

شاخص ۳۰ شرکت بزرگ در ۳ روز گذشته و بازده شاخص ۵۰ شرکت فعال بورس در ۲ روز گذشته، تأثیر قابل توجهی در بهبود عملکرد مدل SVR داشته‌اند. نتایج در جدول ۶ ارائه شده است. با مقایسه نتایج حاصل از MID و FCD، می‌توان گفت MID نتیجه بهتری در الویت‌بندی ویژگی‌ها ارائه نموده است. براین اساس، الگوریتم MID برای اولویت‌بندی و انتخاب ویژگی‌های مناسب برای مدل SVR به‌منظور پیش‌بینی شاخص بورس انتخاب گردید.

**جدول ۵. پانزده ویژگی ابتدایی الویت‌بندی شده توسط MID**

شماره ویژگی	۱	۲	۳	۴	۵
نام ویژگی	تفاوت نسبی درصد	تغییر نسبی حجم	نسبت قیمت به سود	تغییر نسبی در قیمت	بازده شاخص ۳۰
	بازده شاخص در ۵ روز گذشته	معاملات در ۱ روز گذشته	هرسهم در ۳ روز گذشته	دلار در ۱ روز گذشته	شرکت بزرگ در ۱ روز گذشته
	(RDP5)	(TV <sub>1</sub> )	(P/E <sub>43</sub> )	(D <sub>1</sub> )	(LCI30 <sub>1</sub> )
شماره ویژگی	۶	۷	۸	۹	۱۰
نام ویژگی	تفاوت نسبی درصد	تغییر نسبی در	تغییر نسبی در قیمت	بازده شاخص ۳۰	تغییر نسبی در قیمت
	بازده شاخص در ۱۰ روز گذشته	قیمت سکه در ۲ روز گذشته	نفت در ۲ روز گذشته	شرکت بزرگ در ۳ روز گذشته	سکه در ۳ روز گذشته
	(RDP10)	(GP <sub>2</sub> )	(OP <sub>2</sub> )	(LCI30 <sub>3</sub> )	(GP <sub>3</sub> )
شماره ویژگی	۱۱	۱۲	۱۳	۱۴	۱۵
نام ویژگی	بازده شاخص مالی در ۱ روز گذشته	تغییر نسبی در	تفاوت نسبی درصد	تغییر نسبی حجم	بازده شاخص ۵۰
	(FI <sub>1</sub> )	قیمت سکه در ۱ روز گذشته	بازده شاخص در ۱۵ روز گذشته	معاملات در ۲ روز گذشته	شرکت فعال بورس در ۲ روز گذشته
	(GP <sub>1</sub> )	(RDP15)	(TV <sub>2</sub> )	(IC50 <sub>2</sub> )	

**جدول ۶. اثر دو ویژگی ۹ و ۱۵ بر عملکرد مدل SVR با ویژگی‌های اولویت‌بندی شده توسط الگوریتم MID**

شماره ویژگی	نام ویژگی	درصد خطا	MAE	MSE	RMSE
شماره ۹	(LCI30 <sub>3</sub> )	۳/۶۷	۰/۰۰۳۴۲۲	۰/۰۰۰۰۳۰	۰/۰۰۵۴۸۷
شماره ۱۵	(IC50 <sub>2</sub> )	۳/۵۵	۰/۰۰۳۳۰۱	۰/۰۰۰۰۲۹	۰/۰۰۵۳۳۹

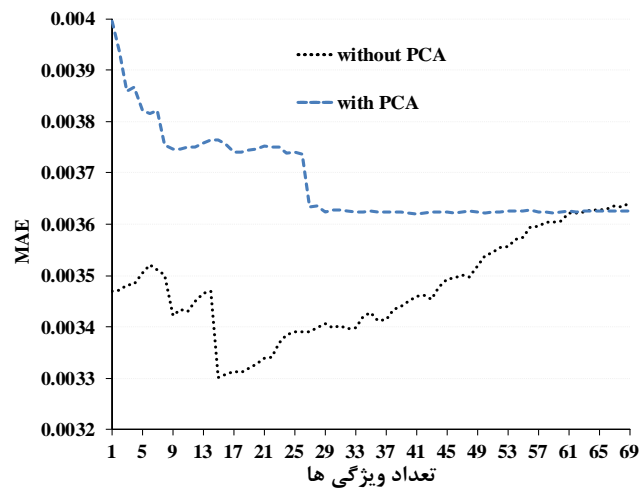
### استخراج ویژگی‌های مؤثر بر شاخص با استفاده از تجزیه و تحلیل مولفه‌های اصلی (PCA)

در این مطالعه، به‌منظور استخراج ویژگی‌های جدید از PCA استفاده می‌شود. در این راستا، ابتدا همه ویژگی‌ها به الگوریتم PCA داده شده و با توجه به تعداد ابعاد خروجی که برای PCA در نظر گرفته می‌شود استخراج ویژگی‌ها انجام می‌شود. به عنوان مثال، اگر دو بعد به عنوان خروجی‌های PCA تعریف گردد، همه‌ی ویژگی‌ها را دریافت کرده و با استفاده از الگوریتم مربوطه دو ویژگی استخراج می‌کند. برای این‌که بتوان راحت‌تر مقایسه استفاده و عدم استفاده از PCA را مشاهده کرد، نتیجه بهره‌گیری از هر دو تکنیک کاهش ابعاد، در نمودار ۲ نمایش داده شده است. با توجه به نتایج مشاهده می‌شود که با افزایش

تعداد ابعاد PCA تا ۲۹ بعد، دقت مدل پیش‌بینی افزایش می‌یابد و بعد از آن تغییر تعداد ابعاد ویژگی‌ها بهبودی در دقت مدل ایجاد نمی‌کند. نتایج در جدول ۷ ارائه شده است. با مقایسه دو منحنی در نمودار ۲، دقت بالاتر استفاده از روش انتخاب ویژگی‌ها نسبت به استخراج ویژگی‌ها کاملاً مشهود است. بنابراین می‌توان به این نتیجه رسید که بدون استفاده از PCA، عملکرد مدل SVR بهتر می‌باشد. بر این اساس، نیاز به استخراج ویژگی برای پیش‌بینی شاخص بورس نیست و باید تمرکز بر روی انتخاب ویژگی‌های مناسب برای آن باشد.

جدول ۷. نتایج استخراج ویژگی با PCA در مدل SVR

RMSE	MSE	MAE	درصد خطا	۲۹ بعد ویژگی
۰/۰۰۵۹	۰/۰۰۰۰۳۶	۰/۰۰۳۶	۳/۸۹	مدل SVR

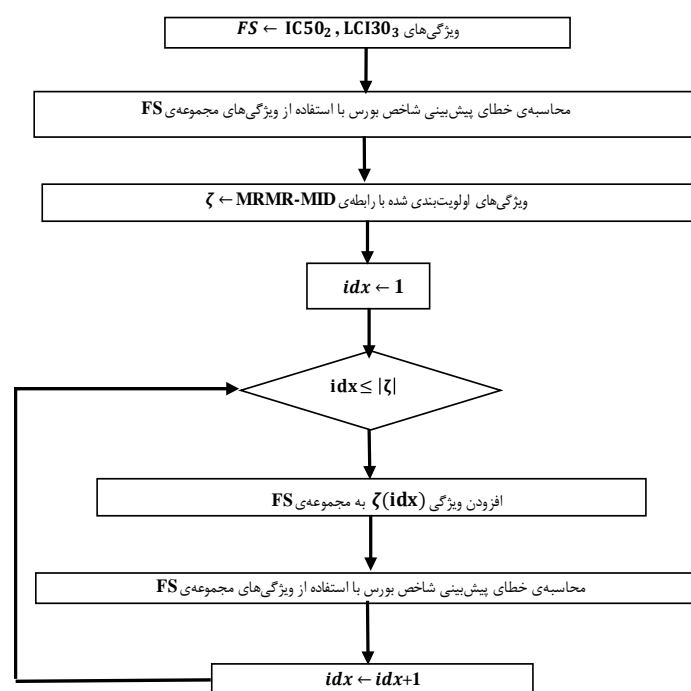


نمودار ۲: استفاده از دو تکنیک کاهش ابعاد برای پیش‌بینی شاخص بورس با بهره‌گیری از مدل SVR

#### الگوریتم پیشنهادی برای انتخاب ویژگی‌های مناسب برای مدل پیش‌بینی شاخص روزانه بورس

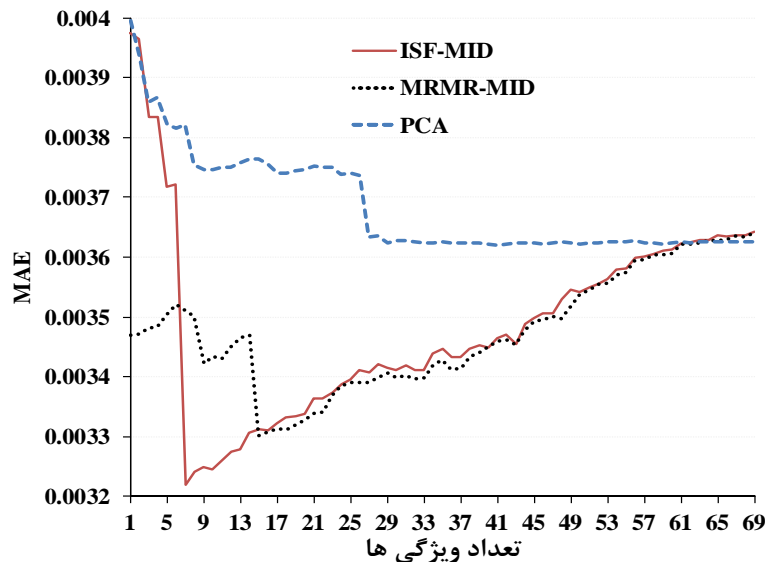
در الگوریتم پیشنهادی برای انتخاب ویژگی‌های مناسب، تعدادی از ویژگی‌ها توسط محقق انتخاب می‌شوند. براساس الگوی تغییر مقدار خطا با افزایش تعداد ویژگی‌ها در روش انتخاب ویژگی با الگوریتم MID، که در نمودار ۱ نمایش داده شده است، مشاهده می‌شود که در مدل SVR، دو ویژگی بازده شاخص ۳۰ شرکت بزرگ در ۳ روز گذشته (ویژگی شماره ۹) و بازده شاخص ۵۰ شرکت فعال بورس در ۲ روز گذشته (ویژگی شماره ۱۵)، باعث کاهش قابل توجه خطای پیش‌بینی شده‌اند. بر این اساس، تصمیم بر آن شد که

برای اولویت‌بندی ویژگی‌ها ابتدا، این دو ویژگی را به‌عنوان ویژگی‌های معرفی شده، انتخاب کرده و سپس به اولویت‌بندی سایر ویژگی‌ها با الگوریتم MID پرداخته شود. ویژگی‌ها با اولویت جدیدی که برای آن‌ها به‌دست آمده است، به مدل SVR به‌منظور پیش‌بینی شاخص بورس به ترتیب داده می‌شوند. فلوجارت الگوریتم پیشنهادی تحت عنوان 'ISF\_MID' در شکل ۱ نمایش داده شده است. شایان ذکر است که در این ارزیابی نیز از روش اعتبارسنجی k-fold (K=10) استفاده می‌شود.



شکل ۱. فلوجارت الگوریتم پیشنهادی برای انتخاب ویژگی‌های مناسب به‌منظور پیش‌بینی شاخص بورس

همچنین نتیجه این محاسبات که اثر افزایش تعداد ویژگی‌ها با اولویت مشخص شده را در دقت مدل SVR نشان می‌دهد، در نمودار ۳ نمایش داده شده است. در این نمودار، برای سهولت مقایسه استخراج ویژگی‌ها با الگوریتم PCA، اولویت‌بندی ویژگی‌ها با الگوریتم MID، و الگوریتم پیشنهادی ISF\_MID، بر دقت عملکرد مدل SVR در پیش‌بینی روزانه شاخص کل بورس اوراق بهادار تهران نشان داده شده است.



نمودار ۳. مقایسه استفاده از PCA، MID، و ISF-MID در انتخاب ویژگی‌های مناسب در مدل SVR به‌منظور پیش‌بینی روزانه شاخص کل بورس اوراق بهادار تهران

با مقایسه نمودارهای موجود می‌توان به این نتیجه رسید که مدل SVR با ویژگی‌های انتخابی توسط الگوریتم ISF\_MID، تنها با ۷ ویژگی، بهترین نتیجه و دقت را ارائه کرده است. در صورتی که مدل SVR با ویژگی‌های انتخابی توسط الگوریتم MID، با ۱۵ ویژگی، و با ویژگی‌های استخراج شده توسط الگوریتم PCA با ۲۹ بعد، بیشترین دقت را دارد. نتایج استفاده از ویژگی‌های انتخابی توسط الگوریتم ISF\_MID، در جدول ۸ نشان داده شده است. با توجه به نتایج ارائه شده در این جدول می‌توان به این نتیجه رسید که؛ مدل SVR با الگوریتم ISF\_MID، بهترین عملکرد را دارد. همان‌طور که ملاحظه می‌شود، مقدار MAE مدل SVR با ویژگی‌های انتخابی توسط الگوریتم ISF\_MID، با ۷ ویژگی ابتدایی، برابر با ۰/۰۰۳ است.

جدول ۸. نتایج انتخاب ویژگی با ISF\_MID در مدل SVR

RMSE	MSE	MAE	درصد خطا	۷ ویژگی
۰/۰۰۵۲	۰/۰۰۰۰۲۸	۰/۰۰۳۲	۳/۴۶	مدل SVR

این ویژگی‌ها به ترتیب عبارتند از: بازده شاخص ۵۰ شرکت فعال بورس در ۲ روز گذشته، بازده شاخص ۳۰ شرکت بزرگ در ۳ روز گذشته، بازده شاخص کل در ۳ روز گذشته، تغییر نسبی حجم معاملات در ۳ روز

گذشته، تفاوت نسبی درصد بازده شاخص در ۱۰ روز گذشته، تغییر قیمت سکه در ۲ روز گذشته، و تفاوت نسبی درصد بازده شاخص در ۵ روز گذشته.

### نتیجه‌گیری و بحث

در این پژوهش سعی شده است با رویکردی کاملاً جدید و منحصر به فرد، ویژگی‌های مناسب برای مدل پیش‌بینی روزانه شاخص کل بورس اوراق بهادار تهران شناسایی شوند. در این راستا به منظور انتخاب ویژگی‌های مناسب برای ورودی مدل SVR، الگوریتمی به نام ISF\_MID، پیشنهاد شد که باعث افزایش دقت مدل ارائه شده به میزان ۹۶/۶۴ درصد گردید. برای شناسایی این الگوریتم ابتدا از هر دو تکنیک کاهش ابعاد (انتخاب و استخراج ویژگی) برای داده‌کاوی و پیش‌پردازش داده‌های ورودی به مدل SVR استفاده گردید، تا تاثیر هر دو روش کاهش ابعاد در دقت عملکرد مدل SVR در پیش‌بینی شاخص روزانه بورس اوراق بهادار تهران شناسایی گردد. برای استخراج ویژگی‌ها از PCA، و برای انتخاب ویژگی‌ها از دو روش تخمین mRMR با نام‌های MID و FCD، استفاده گردید. با مشخص شدن الویت ویژگی‌های انتخاب شده توسط MID و FCD، به مقایسه میزان دقت مدل SVR در استفاده از هر دو الگوریتم انتخاب ویژگی پرداخته شد، و براساس نتایج، MID به عنوان الگوریتم بهتر در انتخاب ویژگی تأیید شد. سپس با مقایسه الگوریتم MID با الگوریتم PCA، عملکرد بهتر استفاده از روش انتخاب ویژگی‌ها نسبت به استخراج ویژگی - ها در انتخاب متغیرهای ورودی مدل SVR مشاهده شد. بر این اساس، نتیجه گرفته شد که ضرورتی به استخراج ویژگی برای پیش‌بینی شاخص بورس نمی‌باشد و باید تمرکز بر روی انتخاب ویژگی‌های مناسب برای آن باشد. در نهایت با توجه به نتایج استفاده از MID در الویت‌بندی ویژگی‌های موثر بر پیش‌بینی شاخص بورس، الگوریتم ISF\_MID، پیشنهاد گردید. با استفاده از این الگوریتم تنها با ۷ ویژگی بهترین نتیجه و دقت را می‌توان با مدل SVR در پیش‌بینی روزانه شاخص کل بورس اوراق بهادار تهران، بدست آورد.

به طور خلاصه؛ اقداماتی که در این پژوهش صورت گرفت که باعث تمایز آن با سایر پژوهش‌ها در این حوزه می‌شود، به شرح زیر می‌باشند:

- جمع‌آوری یک مجموعه کامل از ویژگی‌های تأثیرگذار بر شاخص بورس براساس مبانی تئوریک و بررسی پژوهش‌های پیشین؛ به طوری که ۶۹ ویژگی از ۱۶ شاخص مورد بررسی قرار گرفت.
- این مقاله، از روش اعتبارسنجی متقابل k-fold، برای انتخاب مجموعه‌های آموزش و آزمایش بهره برده‌است. این روش باعث می‌شود همه داده‌ها یک بار برای آموزش و یک بار برای آزمایش به کار روند، که این امر منجر به افزایش دقت و مفید واقع شدن مدل مورد نظر در عمل می‌گردد.
- مطالعات موجود در حوزه پیش‌پردازش داده‌ها به منظور انتخاب ویژگی‌های مناسب برای ورودی مدل پیش‌بینی، اتکا به استفاده از یک نوع تکنیک کاهش ابعاد (انتخاب یا استخراج ویژگی) را نشان می‌دهند. این امر ممکن است برخی از مفروضات مهم در مورد عملکرد رگرسیون اصلی متصل به متغیرهای ورودی و خروجی را نادیده بگیرد. بنابراین در این مطالعه، همزمان از دو

شیوه متفاوت کاهش ابعاد استفاده می‌شود. به طوری که برای الویت‌بندی و انتخاب ویژگی‌های مناسب از mRMR-MID و mRMR-FCD استفاده می‌شود، و برای استخراج ویژگی‌ها PCA به کار گرفته می‌شود. بنابراین برای نخستین بار این پژوهش یک فرآیند داده‌کاوی جامع را برای پیش‌بینی شاخص روزانه بورس انجام می‌دهد.

- برخلاف پژوهش‌های صورت گرفته در این حوزه که از مقدار همبستگی ویژگی‌های ورودی نسبت به خروجی مدل برای انتخاب ویژگی‌های مؤثر استفاده می‌شود، در این پژوهش برای انتخاب ویژگی‌های مناسب از mRMR استفاده گردید. این روش آماری؛ ویژگی‌های مؤثر با توجه به بیشینه‌سازی معیار وابستگی آماری مجموعه ویژگی‌ها با ویژگی هدف، و کمینه کردن اطلاعات متقابل در بین مجموعه ویژگی‌های انتخابی، گزینش می‌شوند.

- الگوریتمی تحت عنوان ISF\_MID برای انتخاب ویژگی‌های مناسب مدل پیش‌بینی شاخص روزانه بورس اوراق بهادار تهران پیشنهاد گردید که با توجه به این الگوریتم می‌توان با تعداد محدودی ویژگی (۷ ویژگی) به پیش‌بینی با بیشترین دقت اقدام نمود، به طوری که میزان درصد خطا در مدل پیشنهادی به ۳/۴۶ رسیده است.

با توجه به نتایج حاصل از پژوهش، پیشنهاد می‌شود به منظور بررسی عملکرد الگوریتم پیشنهادی ISF\_MID، مجموعه داده‌های مورد بررسی را گسترش داد و/ یا از دیگر تکنیک‌های هوشمند به عنوان مدل پیش‌بینی استفاده نمود، و مجدد عملکرد الگوریتم پیشنهادی در این پژوهش را مورد آزمون قرار داد. همچنین برای تکمیل آن توصیه می‌شود در پژوهش‌های بعدی از داده‌های بازار بورس دیگر کشورها استفاده نموده و این الگوریتم در این بازارها مورد مطالعه قرار گیرد. در پایان خاطر نشان می‌کنیم که از نتایج این مطالعه پژوهشگران حوزه مدل‌سازی و سرمایه‌گذاران اعم از حقیقی و حقوقی می‌توانند استفاده کنند.

### ملاحظات اخلاقی

حامی مالی: مقاله حامی مالی ندارد.

مشارکت نویسندگان: تمام نویسندگان در آماده‌سازی مقاله مشارکت داشته‌اند.

تعارض منافع: بنا بر اظهار نویسندگان در این مقاله هیچ‌گونه تعارض منافی وجود ندارد.

تعهد کپی‌رایت: طبق تعهد نویسندگان حق کپی‌رایت رعایت شده است.

## منابع

- Bajalan, S., Fallahpour, S., Dana, N. (2017). "Predicting stock price trends using a modified support vector machine with hybrid feature selection". *Financial Management Perspective*, 7(17), 69-86. (In Persian).
- Bustos, O. Pomares-Quimbaya, A. (2020). "Stock Market Movement Forecast: A Systematic Review", *Expert Systems with Applications*, Volume 156, 15 October, 113464.
- Cavalcante, R. C., Brasileiro, R. C. , Souza V. L.F., Nobrega, J. P. & Oliveira A. L.I. (2016). "Computational Intelligence and Financial Markets: A Survey and Future Directions", *Expert Systems with Applications*. 55.194-211.
- Ding, C. and H. Peng (2005). "Minimum redundancy feature selection from microarray gene expression data". *Journal of bioinformatics and computational biology*. 3(2), 185-205.
- Guo-Qiang, X. (2011). "The optimization of share price prediction model based on 1712 support vector machine". In *International conference on control, automation and 1713 systems engineering* (pp. 1-4). IEEE.
- Henrique, B. M., Sobreiro, V. A., & Kimura, H. (2019). "Literature review: Machine learning techniques applied to financial market prediction". *Expert Systems with Applications*. Volume 124, 15 June. 226-251.
- Huang, C.-F. (2012). "A hybrid stock selection model using genetic algorithms and support vector regression". *Applied Soft Computing*, 12 (2), 807-818.
- Kara, Y. , Boyacioglu, M. A. , & Baykan, O. K. (2011). "Predicting direction of stock price index movement using artificial neural networks and support vector machines: The sample of the Istanbul stock exchange". *Expert Systems with Applications*, 38 (5), 5311-5319.
- Kumar, Deepak. Sarangi, Pradeepta Kumar & Verma, Rajit. (2021). "A systematic review of stock market prediction using machine learning and statistical techniques", *Materials Today: Proceedings*.
- Lee, Ming.Chi (2009). "Using support vector machine with a hybrid feature selection method to the stock trend prediction". *Expert Systems with Applications*. Volume 36. Issue 8, 10896-10904.
- Lui, Y., and Zheng, Y.F. (2006). "FS\_SFS: A novel feature selection method for support vector machines". *Pattern Recognition*. Volume 39, Issue 7, July 2006, Pages 1333-1345.
- Mandal. M and Mukhopadhyay. A. (2013). "An improved minimum redundancy maximum relevance approach for feature selection in gene expression data". *Procedia Technol.*10, 20-27.
- Mansourfar, Gholamreza. Ghayour, Farzad, Khaleghparast Athari, Shabnam. (2015). "Predicting the Industry Index Volatility of Companies Listed in Tehran Stock Exchange, Emphasizing on Corporate Financial Variables Using Support Vector



Machine". *Journal of Empirical Studies in Financial Accounting*, Volume: 12 Issue: 46. (In Persian).

Monajemi, Amirhassan Ebrazi, Medi & Rayati, Alireza. (2009). "Stock price prediction in Tehran stock exchange using artificial neural network". *Journal of financial economy*, 6(3), 1-26. (In Persian).

Nevasalmi, Lauri. (2020). "Forecasting multinomial stock returns using machine learning methods". *The Journal of Finance and Data Science*, Volume 6, 86-106.

Nguyen, Duc-Hien, Le Manh-Thanh. (2014). "A two-stage architecture for stock price forecasting by combining SOM and fuzzy-SVM", *International Journal of Computer Science and Information Security (IJCSIS)*, Vol. 12, No. 8, August.

Ni, L.P., Ni, Zh. W., & Gao, Y.Zh. (2011). Stock trend prediction based on fractal feature selection and support vector machine. *Expert Systems with Applications*, 38(5): 5569-5576.

Ou, P., & Wang, H. (2009). "Prediction of stock market index movement by ten data mining techniques". *Modern Applied Science*, 3, P28.

Patel, J., Shah, S., Thakkar, P., and Kotecha, K. (2015). "Predicting stock and stock price index movement using trend deterministic data preparation and machine learning techniques". *Expert Systems with Applications*, 42(1):259-268.

Pearson, K. (1901). "On lines and planes of closest fit to systems of points in space". *Philosophical Magazine*, 2(6), 559-572.

Perez-Rodriguez, J. V., S. Torrab and J. Andrada-Felixa (2004). "STAR and ANN models: Forecasting performance on the Spanish Ibex-35 stock index". *Journal of Empirical Finance*. 12(3), 490-509.

Raee, R., Nikahd, A., Habibi, M. (2017). "The Index Prediction of Tehran Stock Exchange by Combining the Principal Components Analysis, Support Vector Regression and Particle Swarm Optimization". *Financial Management Strategy*, 4(4), 1-23. (In Persian).

Rafiuzzaman, M. (2014). "Forecasting Chaotic Stock Market Data using Time Series Data Mining". *International Journal of Computer Applications*. 101(10), 27-34.

Singh, R. and Srivastava, S. (2017). "Stock prediction using deep learning". *Multimedia Tools and Applications*, 76(18):18569-18584.

Ul Haq, Anwar. Zeb, Adnan. Lei, Zhenfeng & Zhang, Defu. (2021). "Forecasting daily stock trend using multi-filter feature selection and deep learning", *Expert Systems with Applications*, 168 (2021) 114444

Wanga, Diya & Zhao, Yixi (2020) "Using News to Predicton Investor Sentiment: Based on SVM Model", *Procedia Computer Science*. Wolume 174 .191-199

Wei, Z. (2012). *A svm approach in forecasting the moving direction Chinese stock indices*, Department of industrial and systems engineering, Thesis of Master of Sciences, Lehigh University.

Yuan, Y. (2013). "Forecasting the movement direction of exchange rate with polynomial smooth support vector machine". *Mathematical and Computer Modelling*, 57 (3), 932–944.

Zhang, X., Hu, Y., Xie, K., Wang, S., Ngai, E. W. T., & Liu, M. (2014). "A causal feature selection algorithm for stock prediction modeling". *Neurocomputing*, 142. 48-59.

Zhong, X., & Enke, D. (2017). "Forecasting daily stock market return using dimensionality reduction". *Expert Systems with Applications*, 67, 126–13.

#### COPYRIGHTS



©2022 Alzahra University, Tehran, Iran. This license allows others to download the works and share them with others as long as they credit them, but they can't change them in any way or use them commercially.