

طراحی مدل پیش بینی تقاضا در صنعت کاشی و سرامیک

علیرضا موتمنی*، مصطفی رضائی**، مریم احقائی***

چکیده

پیش‌بینی تقاضا جزء مهم‌ترین فعالیتهای هر سازمان برای برنامه‌ریزی فروش و در نهایت برنامه‌ریزی جامع بوده و در واقع تعیین‌کننده حجم فعالیتهای سازمان در آینده می‌باشد. همچنین درک درستی از میزان و کیفیت فعالیتهای مزبور را برای مدیران فراهم می‌نماید. در این مقاله بر اساس مبانی علوم هوش مصنوعی و داده‌کاوی، به ارائه مدلی جهت پیش‌بینی میزان فروش در صنعت کاشی و سرامیک پرداخته شده است. مدل پیشنهادی، یک مدل ترکیبی شامل کاهش بعد، خوشه‌بندی و پیش‌بینی بوده و به منظور اجرای مراحل مختلف آن از الگوریتم‌های آنالیز مؤلفه‌های مستقل، یادگیری منیفلد، خوشه‌بندی کامینز و رگرسیون بردار پشتیبان، استفاده شده است. در تحقیق حاضر ۵۰ مورد از فروش‌های ماهیانه مربوط به ۳ سال گذشته شرکت کاشی ایرانا مورد استفاده قرار گرفته است. نتایج بدست آمده از طریق پیش‌بینی فروش با مدل پیشنهادی، بدلیل کاهش خطاهای عمومی و خطاهای نمونه در مقایسه با مدل‌های سنتی پیش‌بینی، از دقت بالاتری برخوردار می‌باشد.

کلیدواژه‌ها: پیش‌بینی تقاضا؛ رگرسیون بردار پشتیبان؛ یادگیری منیفلد؛ آنالیز مؤلفه‌های مستقل؛ خوشه‌بندی کامینز.

تاریخ دریافت مقاله: ۹۱/۱۲/۲۴، تاریخ پذیرش مقاله: ۹۲/۴/۳۱.

* استادیار، دانشگاه شهید بهشتی (نویسنده مسئول).

E-mail: ar_motameni@yahoo.com

** کارشناسی ارشد، دانشگاه سمنان.

*** کارشناسی ارشد، دانشگاه تربیت مدرس.

۱. مقدمه

یکی از مسائل چالش‌برانگیز در بهبود عملکرد سازمان‌ها، پیش‌بینی تقاضا، بهبود زنجیره تأمین آن‌ها و کاهش هزینه‌های مرتبط با آن است. بنا بر آمار منتشر شده، در سال ۱۹۹۸ شرکت‌های آمریکایی، ۸۹۸ میلیارد دلار صرف هزینه‌های مرتبط با زنجیره تأمین کرده‌اند (این مبلغ با ۱۰ درصد تولید ناخالص ملی آن کشور در آن سال برابر است). این عدد به یک تریلیون دلار در سال ۲۰۰۰ افزایش یافت و در سال‌های ۲۰۰۱ و ۲۰۰۲ برابر با ۹۵۷ میلیارد دلار و ۹۱۰ میلیارد دلار گزارش شده است [۱]. متأسفانه بیشتر این هزینه‌ها غیرضروری بوده و با بهبود مدیریت زنجیره تأمین می‌توان به میزان قابل توجهی هزینه‌های مزبور را کاهش داد [۲:۷]؛ برای نمونه متخصصین معتقدند که با بهره‌گیری از استراتژی‌های بهتر زنجیره تأمین، صنایع غذایی می‌توانند حدود ۳۰ میلیارد دلار یا ۱۰ درصد از هزینه‌های عملیاتی سالانه خود بکاهند [۳]. از جمله راه‌حل‌های کاهش هزینه‌های زنجیره تأمین، افزایش دقت پیش‌بینی تقاضا به کمک روش‌های جدید آماری است. روش‌های سنتی پیش‌بینی تقاضا علاوه بر دقت پایین، مشکلات دیگری نیز برای سازمان‌ها به وجود می‌آورند؛ از جمله این مشکلات، می‌توان از «اثر شلاقی» یاد کرد [۲:۲۳]. با پیشرفت‌های اخیر در هوش مصنوعی، تکنیک‌های جدیدی برای پیش‌بینی ارائه شده است که نسبت به تکنیک‌های سنتی، از دقت بالاتری برخوردار هستند. رایج‌ترین این تکنیک‌ها، الگوریتم‌های شبکه عصبی است که کاستی‌هایی مانند نیاز به پارامترهای کنترلی زیاد، دشواری رسیدن به نتیجه‌ای پایدار و خطر برازش بیش از حد، دارد [۴، ۵، ۶]. به دلیل وجود چنین ضعف‌هایی، مدل‌های بهتری برای بهبود مدل شبکه عصبی طراحی شده است. ماشین‌های بردار پشتیبان از الگوریتم شبکه عصبی بدیعی بر مبنای تئوری یادگیری آماری بهره می‌گیرند [۷، ۸]. این الگوریتم‌ها پتانسیل بسیار بالایی داشته و در مسائل کاربردی نتایج بسیار خوبی به دنبال دارند. امروزه مدل رگرسیون ماشین‌های بردار پشتیبان به نام «رگرسیون بردار پشتیبان» نیز در موضوع پیش‌بینی‌های غیرخطی، بسیار مورد توجه قرار گرفته و نتایج بسیار خوبی به همراه داشته است. با این حال بسیاری از محققان یادآور این نکته شده‌اند که وجود اغتشاش در داده‌ها، نتایج پیش‌بینی را تا حد بسیار زیادی تحت تأثیر قرار می‌دهد [۹، ۱۰]. به نظر می‌رسد پیش‌پردازش داده‌ها بیش از هر زمانی اهمیت پیدا کرده است [۱۱].

با توجه به اهمیت موضوع پیش‌بینی تقاضا در سازمان‌ها، در این مقاله مدل سه‌مرحله‌ای برای پیش‌بینی فروش ارائه شده است به این صورت که قبل از انجام پیش‌بینی بوسیله رگرسیون بردار پشتیبان از دو مرحله خاص برای پیش‌بینی پردازش داده‌ها استفاده گردیده است. در مرحله اول سعی می‌شود که با کاهش بعد داده‌ها اغتشاشات موجود در آن‌ها تا حد امکان گرفته شود و، در مرحله دوم با دسته‌بندی و خوشه‌بندی کردن داده‌ها و قرار دادن داده‌های مشابه

در گروه‌های یکسان، دقت پیش‌بینی افزایش داده شده و در نهایت پیش‌بینی را روی داده‌هایی که اغتشاشات آنها تا حد امکان گرفته شده انجام می‌گیرد.

۲. پیشینه تحقیق

در انجام این تحقیق از مفاهیم آنالیز اجزای اصلی، یادگیری منیفلد، آنالیز اجزای مستقل، خوشه‌بندی کامینز، ماشین بردار پشتیبان، رگرسیون بردار پشتیبان و بیش‌برازش استفاده شده که در ادامه، تعریف مختصری از مفاهیم مزبور و نتایج تحقیقات مرتبط ارائه گردیده است.

آنالیز اجزای اصلی^۱. آنالیز اجزای اصلی از روش‌های کلاسیک در آنالیز آماری داده^۲، استخراج ویژگی^۳ و فشرده‌سازی داده به‌شمار می‌رود که از لحاظ تاریخی، به کارهای اولیه پیرسون^۴ در حدود سال‌های ۱۹۰۰ بازمی‌گردد.

هدف اصلی آنالیز اجزای اصلی، پیدا کردن راستاهای مؤلفه‌های اصلی است. راستاهای مؤلفه‌های اصلی راستاهایی هستند که تصویر داده‌ها بر آن‌ها حداکثر واریانس را خواهند داشت؛ به این ترتیب، مؤلفه اصلی اول، راستایی است که تصویر داده‌ها بر آن حداکثر واریانس را دارد، دومین مؤلفه اصلی، راستایی است که از میان تمامی راستاهای عمود بر راستای اول، تصویر داده‌ها ماکزیمم واریانس را دارا بوده و k امین مؤلفه‌ی اصلی، راستایی است که از میان تمامی راستاهای عمود بر $k-1$ راستای قبل، تصویر داده‌ها بر آن بیشترین واریانس را دارد.

برای پیدا کردن راستای همه مؤلفه‌های اصلی، محاسبات به صورت جبر خطی انجام خواهد شد. اگر تمامی n داده مورد مطالعه را روی هم بگذاریم و در ماتریس X که یک ماتریس $n \times p$ است، بنویسیم، تصویر داده‌ها Xw خواهد بود که یک ماتریس $n \times 1$ بوده و واریانس آن برابر خواهد بود با:

$$\sigma_w^2 = \frac{1}{n} \sum_i (\vec{x}_i \cdot \vec{w})^2 \quad \text{رابطه ۱}$$

-
1. Principle Component Analysis
 2. Statistical Data Analysis
 3. Feature Extraction
 4. Pearson

$$\begin{aligned}
 &= \frac{1}{n} (Xw)^T (Xw) \\
 &= \frac{1}{n} w^T X^T X w \\
 &= w^T \frac{X^T X}{n} w \\
 &= w^T V w
 \end{aligned}$$

می‌توان نشان دهیم که بردار w مورد نظر، بردار ویژه ماتریس کوواریانس V است. از آنجا که V یک ماتریس $P \times P$ است، تعداد P بردار ویژه متفاوت خواهد داشت.

یادگیری منیفلد. منیفلد یک فضای فرعی منحنی و چنبری است که در درون یک فضای اقلیدسی، جاسازی شده است.

نکته مهم در یادگیری منیفلد این است که هر منیفلد q بعدی را می‌توان با یک فضای فرعی خطی q بعدی تخمین زده و در اطراف هر نقطه فضای کوچکی را در نظر گرفته و فضای مماس بر آن را پیدا نمود. همچنین فضای مماس نقاط در طول منیفلد به صورت پیوسته تغییر می‌کند. هر چه تغییرات فضای مماسی بیشتر باشد، انحنای منیفلد نیز بیشتر خواهد بود. بنابراین اگر داده‌های ما روی یک منیفلد قرار گرفته باشد، می‌توانیم تخمین خطی از منیفلد به صورت محلی داشته باشیم.

بایستی توجه داشت که در یادگیری منیفلد برای کاهش بعد غیر خطی از روشهای مختلفی استفاده می‌شود. یکی از روش‌های کاهش بعد غیرخطی الگوریتم 1 LLE است که به اختصار در زمینه روش مزبور توضیح داده می‌شود.

الگوریتم LLE سه مرحله دارد؛ در مرحله اول برای هر یک از نقاط یک همسایگی تشکیل داده می‌شود، در مرحله دوم برای تخمین خطی نقاط بر اساس همسایگی آن‌ها اقدام گردیده و محاسبه وزن‌ها تخمین زده شده و در مرحله سوم بر اساس وزن‌های پیدا شده، مختصات منیفلد پیدا می‌گردد.

مراحل الگوریتم LLE به صورت زیر است:

۱. برای هر داده k نزدیک‌ترین همسایه شاخص می‌شود.

1. Locally Linear Embedding

۲. ماتریس وزن w به صورتی پیدا می گردد که جمع مجذور خطاهای بازسازی داده ها از همسایه هایش حداقل شود:

رابطه ۲

$$RSS(w) \equiv \sum_{i=1}^n \left\| \vec{x}_i - \sum_{j \neq i} w_{ij} \vec{x}_j \right\|^2$$

۳. مختصات Y را به صورتی مشخص می شود که بازسازی نقاط توسط وزن ها، حداقل گردد.

آنالیز اجزای مستقل^۱. آنالیز اجزای مستقل یک روش پردازش سیگنال آماری است که برای یافتن منابع مستقل (تنها با در دست داشتن داده هایی که ترکیب منابع ناشناخته هستند و بدون داشتن اطلاعات از مکانیزم ترکیب) طراحی شده است.

در این مدل داده های مشاهده شده به صورت $X=AS$ نمایش داده می شود که در آن A ماتریس ناشناخته ای است که نحوه ی ترکیب را نشان می دهد و ماتریس ترکیب نامیده می شود و S منابعی هستند که به طور مستقیم از روی X قابل مشاهده نیستند.

مدل آنالیز اجزای مستقل نشان می دهد که داده های مشاهده شده X از ترکیب خطی منابع مکنون S حاصل شده و ماتریس ترکیب A چگونگی ترکیب خطی را نشان می دهد. فرض اصلی این مدل این است که منابع مکنون از نظر آماری مستقل از یکدیگرند. بر پایه این فرض طی یک فرایند یادگیری بدون نظارت^۲ آنالیز اجزای مستقل به کشف ماتریس تفکیک W می رسد.

بعد از محاسبه ماتریس تفکیک W از آن برای انتقال داده های مشاهده شده X به منابع مستقل Y استفاده می کند ($WX=Y$) و این منابع مستقل Y تخمینی برای منابع مکنون S محسوب می شوند. به ردیف های ماتریس Y اجزای مستقل^۳ گفته می شود و تا جایی که امکان دارد باید نسبت به هم مستقل باشند.

هرچند از آنالیز اجزای مستقل در حل مسائل پردازش سیگنال در امور پزشکی، در زمینه پردازش سیگنال صوتی، تشخیص چهره استفاده زیادی شده است، اما کمتر در پیش بینی تقاضا به کار رفته است.

بک و ویگند [۱۲] از آنالیز اجزای مستقل برای کشف متغیرهای اصلی در تعیین سود روزانه ۲۸ سهام بزرگ ژاپن استفاده کرده اند، کیوی لوتو و آجا [۱۳] این آنالیز را برای کشف عوامل

1. Independent Component Analysis
1. Unsupervised Learning
2. Independent Components

اصلی تعیین‌کننده گردش مالی ۴۰ مغازه در یک مجتمع تجاری به کار برده‌اند و آجا و همکاران [۱۴] از آن برای پیش‌بینی نرخ تبادل ارز استفاده کرده‌اند.

الگوریتم خوشه‌بندی کامینز^۱. الگوریتم خوشه‌بندی کامینز، داده‌ها را بر مبنای ویژگی‌هایشان به K خوشه جداگانه تقسیم می‌کند. اشیائی که در یک خوشه قرار می‌گیرند، ویژگی‌های مشابه دارند. در زیر مراحل الگوریتم خوشه‌بندی کامینز آمده است:

۱. تعداد خوشه‌ها (K) مشخص می‌شود.
 ۲. مرکز ثقل برای خوشه‌ها تعیین می‌گردد.
 ۳. حلقه‌ای روی تمام اشیاء اجرا نموده و فاصله هر شیء را تا مرکز ثقل همه خوشه‌ها اندازه گرفته شده و هر شیء به خوشه نزدیک‌ترین مرکز ثقل نسبت داده می‌شود.
 ۴. مراکز ثقل خوشه‌های جدید مجدداً محاسبه می‌گردد.
 ۵. مرحله ۳ مجدداً اجرا شده تا زمانی که مراکز ثقل دیگر تغییر نکنند.
- در این الگوریتم فاصله بین دو شیء توسط تابع فاصله اقلیدسی (رابطه ۳) محاسبه می‌شود.

$$d(x, y) = \sqrt{\sum_{i=1}^m (x_i - y_i)^2} \quad \text{رابطه ۳}$$

تی و کاو [۵] مدل پیش‌بینی سری زمانی را با ترکیب SVM و خوشه‌بندی ارائه دادند. کاو [۱۵] SVM و خوشه‌بندی را برای مسائل پیش‌بینی سری زمانی به کار برده است. لای و همکاران [۱۶] از الگوریتم K-means به همراه درخت تصمیم‌گیری فازی برای پیش‌بینی قیمت سهام استفاده کرده‌اند. هوانگ و سای [۱۷] از مدل هیبریدی شامل مراحل خوشه‌بندی، SVR و انتخاب متغیر بر مبنای فیلتر برای کاهش زمان برازش مدل و بهبود دقت مدل استفاده کرده‌اند.

ماشین بردار پشتیبان^۲ و رگرسیون بردار پشتیبان^۳. ماشین بردار پشتیبان یکی از روش‌های یادگیری با نظارتی^۴ است که برای طبقه‌بندی و رگرسیون استفاده می‌شود. ماشین بردار پشتیبان در سال ۱۹۹۲ توسط وپنیک^۵ معرفی شده و بر پایه تئوری یادگیری آماری بنا

1. K-Means
 2. Support Vector Machines
 3. Support Vector Regression
 4. Supervised Learning
 5. Vapnik^{ik}

نهاده شده است. ماشین بردار پشتیبان در انواع دسته‌بندی‌ها همچون تشخیص ارقام دستنویس، تشخیص، شناسایی صورت، دسته‌بندی انواع صداها و مانند آن مورد استفاده قرار گرفته است که در مقایسه با تکنیک‌های دیگر از کارایی قابل توجهی برخوردار است. رویکرد SVM به این صورت است که در فاز آموزش، تلاش می‌شود مرز تصمیم‌گیری به گونه‌ای انتخاب شود که حداقل فاصله آن با هر یک از دسته‌های موردنظر حداکثر شود. این نوع انتخاب باعث می‌شود که تصمیم‌گیری در عمل، شرایط نویزی را به خوبی تحمل کرده و پاسخ‌دهی خوبی داشته باشد. این نحوه انتخاب مرز بر اساس نقاطی به نام بردارهای پشتیبان انجام می‌شود. تفاوت اساسی این طبقه‌بندی‌کننده با سایر طبقه‌بندی‌کننده‌های آماری این است که برای پردازش و طبقه‌بندی داده‌های ابرطیفی، دیگر نیازی به کاهش تعداد باندها نمی‌باشد. مسئله رگرسیون غیرخطی می‌تواند مانند مسئله بهینه‌سازی بیان شود.

رابطه ۴

$$\text{Minimize } \frac{1}{2} w^T + C \sum_{i=1}^1 (\xi_i + \xi_i^*)$$

با لحاظ محدودیت‌های زیر

رابطه ۵

$$\begin{aligned} y_i - (w \cdot \varphi(x_i) + b) &\leq \varepsilon + \xi_i \\ (w \cdot \varphi(x_i) + b) - y_i &\leq \varepsilon + \xi_i^* \\ \xi_i, \xi_i^* &\geq 0, i = 1, 2, \dots, 1 \end{aligned}$$

بنابراین مسئله بهینه‌سازی در حالت غیرخطی به صورت زیر خواهد بود:

رابطه ۶

$$\begin{aligned} \text{Minimize } \frac{1}{2} \sum_{ij} (a_i - a_i^*)(a_j - a_j^*) \langle \varphi(x_i) \cdot \varphi(x_j) \rangle \\ + \varepsilon \sum_{i=1}^1 (a_i + a_i^*) - \sum_{i=1}^1 y_i (a_i - a_i^*) \end{aligned}$$

با لحاظ محدودیت‌های زیر

رابطه ۷

$$\sum_{i=1}^l (a_i - a_i^*) = 0$$

$$0 \leq a_i \leq C, i = 1, 2, \dots, l$$

$$0 \leq a_i^* \leq C, i = 1, 2, \dots, l$$

پارامترهای حاکم بر SVR غیرخطی، مقدار ثابت C ، شعاع لوله غیرحساس ε و پارامترهای هسته‌ای هستند. این پارامترها به طور متقابل به هم وابسته بوده و تغییر مقدار از یک پارامتر، به پارامترهای دیگر را تغییر می‌دهد. پارامتر C ، شفافیت تابع تقریب را کنترل می‌کند. مقدار C بزرگ‌تر، خطای بیشتری را ایجاد نموده و یادگیری را پیچیده‌تر می‌کند و یک مقدار C کوچک‌تر به خطاهایی منجر می‌شود که قابل قبول بوده لیکن ممکن است یادگیری با تقویت ضعیف را ایجاد کند.

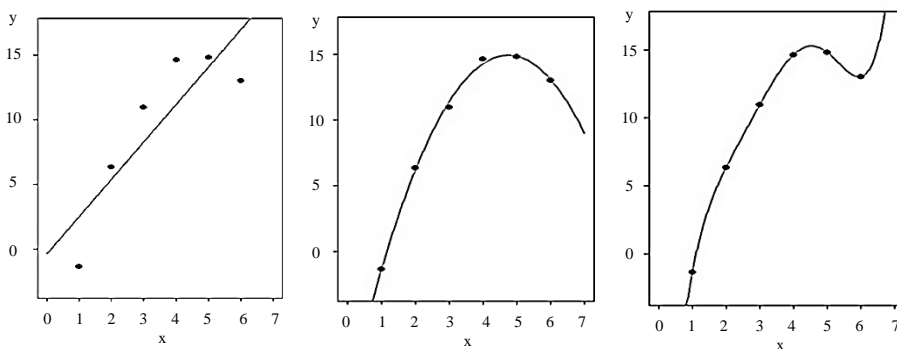
اگر تعداد داده‌ها زیاد شد، ممکن است مقادیر C کوچک‌تر که خطای کمتری را ایجاد می‌کند، ترجیح داده شود. پارامتر ε نیز بر پیچیدگی یا شفافیت تابع تقریب اثر می‌گذارد و حاکم بر تعداد بردارهای پشتیبان است. مقادیر کوچک‌تر ε ممکن است به بردارهای پشتیبان بیشتر و در نتیجه یک ماشین یادگیری پیچیده‌تر منجر شود و مقادیر بزرگ‌تر ε سبب می‌شود که ε -tube تعداد داده بیشتری را احاطه کند که در یادگیری در نظر گرفته نمی‌شود، بنابراین مقدار زیادی از اطلاعات مهم در داده‌ها از دست می‌رود.

رگرسیون بردار پشتیبان به طور موفقیت‌آمیز در مسائل پیش‌بینی سری زمانی و پیش‌بینی تقاضا استفاده شده است؛ برای نمونه در پیش‌بینی ارزش تولید ماشین‌آلات صنعتی، پیش‌بینی ضریب اطمینان موتور، پیش‌بینی سرعت باد و پیش‌بینی سری‌های زمانی مالی [۵، ۱۸، ۱۹، ۲۰، ۲۱، ۲۲، ۲۳ و ۲۴].

اگرچه استفاده از رگرسیون بردار پشتیبان در پیش‌بینی سری‌های زمانی مالی نتایج خوبی به همراه داشته است، تحقیقات کمی در بهره‌گیری از رگرسیون بردار پشتیبان برای پیش‌بینی تقاضا شده است.

بیش برآزش. اگر به مدلی برای پیش‌بینی تقاضا، شامل X مقدار تقاضا و Y مبلغ تقاضا، به صورت شکل ۱ نیاز باشد، برای رسیدن به مدل مناسب پیش‌بینی تقاضا، از این ۶ نقطه رگرسیون گرفته می‌شود:

شکل ۱ نشان‌دهنده نتیجه سه رگرسیون متفاوت است؛ شکل سمت راست یک رگرسیون درجه ۵، شکل وسط رگرسیون درجه ۲ و شکل سمت چپ رگرسیون درجه ۱ را نشان می‌دهد.



شکل ۱. رگرسیون‌های مختلف بر اساس مقدار و مبلغ تقاضا

همان طور که دیده می شود، شکل سمت راست (رگرسیون درجه ۵) از تمامی داده‌ها عبور می کند و بنابراین خطای نمونه این مدل صفر است. اگر چه این مدل، مدل مناسبی برای پیش بینی ۶ مقدار اندازه گرفته شده است، اما مدل مناسبی برای پیش بینی مقادیر خارج نیست؛ به سخن دیگر اگر چه مدل سمت راست خطای نمونه پایینی دارد، اما خطای عمومی آن بالا خواهد بود. مدل سمت چپ مدل بسیار ساده‌ای است که انتظار می رود خطای عمومی آن نیز بالا باشد. گرچه با اینکه خطای عمومی هر دو مدل سمت راست و چپ بالا است، اما خطای بالای آن‌ها به دلایل کاملاً متفاوتی بر می گردد.

در صورتی که رابطه اصلی میان X و Y یک رابطه غیرخطی باشد، حتی اگر داده‌های بسیار زیادی در دست باشد، باز هم مدل سمت چپ (رگرسیون یک جمله‌ای) نخواهد توانست ساختار موجود در داده‌ها را کشف کند.

از طرف دیگر مدل های پیچیده سمت راست (رگرسیون پنج جمله‌ای)، در داده‌های کم، ساختارهایی را که به صورت تصادفی و بر حسب اتفاق در داده‌های ما وجود دارد، به اشتباه به عنوان ساختار عمومی رابطه X و Y بیان خواهد کرد.

خطای بالای عمومی مدل سمت راست به دلیل استفاده کردن از مدلی بسیار پیچیده‌تر از حد مطلوب است و به اصطلاح «بیش‌برازش^۱» رخ داده است.

از سوی دیگر، مدل سمت چپ بسیار ساده است و به همین دلیل خطای عمومی بالایی دارد و به اصطلاح در مدل سمت چپ «کم‌برازش^۲» رخ داده است.

مدل وسط (رگرسیون درجه ۲) از میزان پیچیدگی مناسبی برخوردار است و به طور کلی میزان خطای عمومی به پیچیدگی مدل بستگی دارد.

1. Over fitting
 2. Under fitting

با افزایش پیچیدگی مدل، خطا در نمونه (به‌طور پیوسته) کاهش یافته و مدل‌های پیچیده‌تر، خطا در نمونه کمتری دارند، اما این وضعیت برای خطای عمومی صادق نیست. همچنین با افزایش پیچیدگی تا مقدار خاصی، خطای عمومی بالا خواهد رفت، زیرا با افزایش پیچیدگی به جای کشف روابط حقیقی میان متغیرها، اغتشاش موجود در داده‌ها پردازش می‌شود و مدل بیش‌برازش می‌شود.

همچنین مدل‌هایی که پیچیدگی کمتری از میزان بهینه دارند، مشکل کم‌برازش دارند. لازم به ذکر است که موضوع خطا در این زمینه از اهمیت زیادی برخوردار است. شناخت اجزای اصلی خطای کلی کمک زیادی به حل موضوع می‌نماید. به‌طور کلی خطای کلی یک مدل به سه جزء زیر تقسیم می‌شود:

$$\text{خطای کلی} = (\text{نویز}) + (\text{بایاس}) + (\text{واریانس}) \quad \text{رابطه ۸}$$

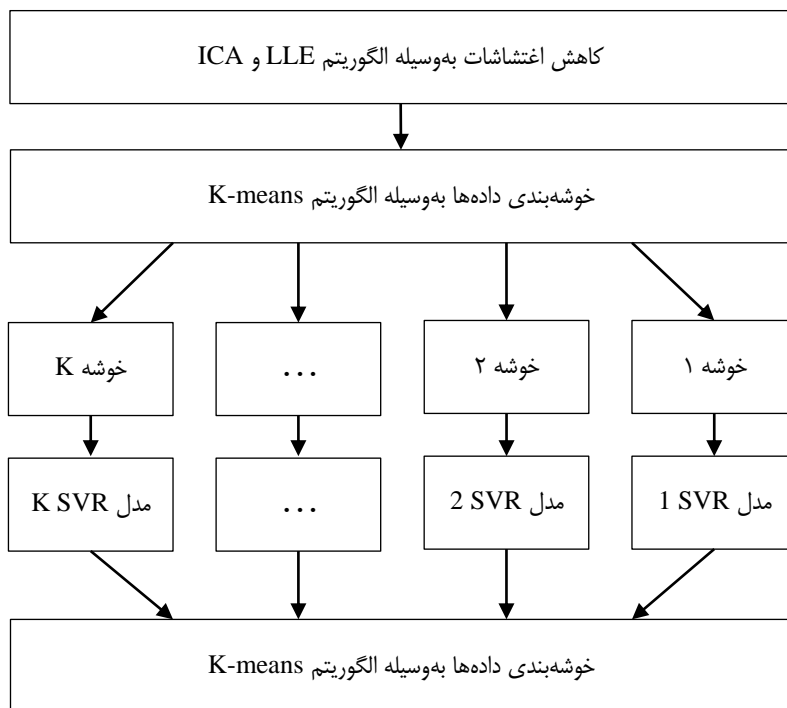
جزء نویز^۱ در اصل جزء غیرقابل پیش‌بینی مسئله محسوب می‌شود و شامل مقدار اغتشاشات آماری^۲ است که حتی در اطراف بهترین پیش‌بینی‌ها وجود دارد، جزء بایاس^۳ باعث می‌شود خطای عمومی در مدل‌هایی که از کم‌برازشی رنج می‌برند، افزایش یابد و جزء واریانس^۴ باعث می‌شود خطای عمومی در مدل‌هایی که از بیش‌برازشی رنج می‌برند، افزایش یابد.

مدل مفهومی. مدل پیشنهادی این مقاله یک مدل ترکیبی^۵ و شامل سه مرحله کاهش بعد، خوشه‌بندی و پیش‌بینی است.

در مرحله اول تلاش شده با کاهش بعد داده‌ها اغتشاش‌های موجود در آن‌ها تا حد امکان گرفته شده و در مرحله دوم با دسته‌بندی و خوشه‌بندی کردن داده‌ها و قرار دادن داده‌های مشابه در گروه‌های یکسان، دقت پیش‌بینی بالا خواهد رفت و در مرحله‌ی سوم پیش‌بینی روی داده‌هایی انجام می‌شود که اغتشاش‌های آن‌ها تا حد امکان گرفته شده است و دسته‌بندی شده‌اند.

-
1. Noise
 2. Statistical Fluctuation
 3. Bias
 4. Variance
 5. Hybrid Model

مراحل اجرای کار در شکل ۲ مشاهده می شود. همان طور که در این شکل مشخص است، بعد از اجرای مراحل سه گانه فوق، به تعداد دسته های به دست آمده از مرحله خوشه بندی، مدل پیش بینی کننده وجود خواهد داشت.

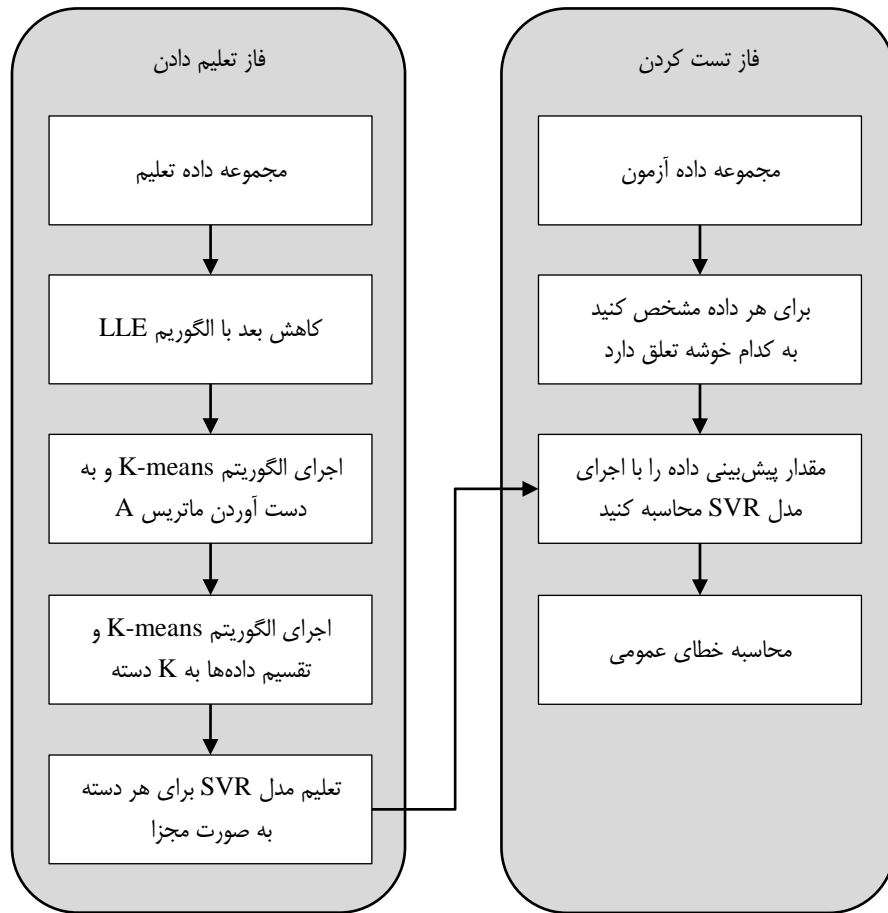


شکل ۲. مراحل اجرای کاهش اغتشاشات

مراحل ذکر شده در بالا باید برای دستیابی به مدل های پیش بینی کننده انجام شوند، اما اگر بخواهیم بعد از به دست آوردن مدل های پیش بینی، میزان فروش داده های جدید را پیش بینی کنیم، باید گام های زیر برداشته شود:

۱. ابتدا باید مشخص شود که داده ای که خواهان پیش بینی آن هستیم به کدام یک از دسته های ۱ تا K تعلق دارد.
۲. بعد از مشخص کردن دسته مورد نظر، از مدل مربوط به آن دسته برای انجام پیش بینی استفاده خواهد شد.

نحوه اجرای پیش بینی در فاز آزمایشی در شکل ۳ نشان داده شده است.



شکل ۳. مراحل تعلیم دادن و آزمایش کردن داده‌ها برای پیش‌بینی

۳. روش‌شناسی

داده‌های به‌کارگرفته شده. شرکت سهامی عام کارخانه چینی ایران (کاشی ایرانا) اولین و بزرگترین تولیدکننده کاشی دیواری به روش صنعتی در ایران است که در سال ۱۳۳۶ با سرمایه ۵۲۰ میلیون ریال تحت شماره ۵۹۰۴ به ثبت رسیده است. این شرکت اکنون با داشتن ۴ سالن تولید، انواع کاشی دیواری به ابعاد ۴۵ در ۳۰، ۶۰ در ۳۰ و ۴۰ در ۲۵ و کاشی کف به ابعاد ۵۰ در ۵۰، ۴۰ در ۴۰، ۳۳ در ۶۰ و ۲۵ در ۵۰ را تولید می‌کند. در تحقیق حاضر از اطلاعات واحد فروش ۵۰ مورد فروش، در طی ۳ سال و به تفکیک ماه استفاده شده است.

این اطلاعات را می توان به صورت یک ماتریس در نظر گرفت؛ به این صورت که هر سطر از این ماتریس، اطلاعات فروش یک مدل و هر ستون، اطلاعات فروش یک ماه را نشان می دهد. از آنجا که در این پژوهش ۵۰ مورد بررسی شده است، لذا ۵۰ سطر خواهیم داشت و از آنجا که از اطلاعات ۳۶ ماه استفاده شده است، ۳۶ ستون خواهیم داشت؛ بنابراین ماتریس به دست آمده، یک ماتریس ۳۶ در ۵۰ خواهد بود.

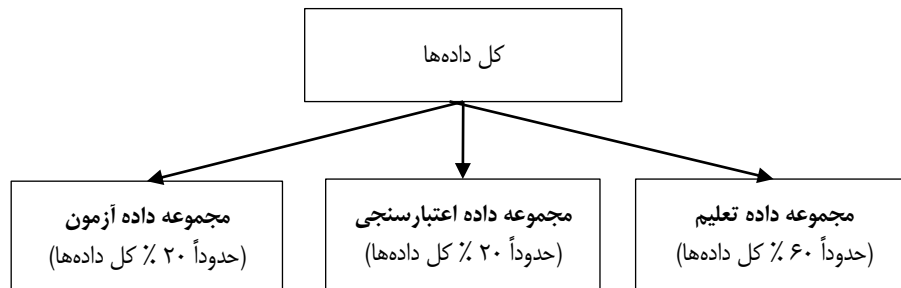
با توجه به اینکه این ماتریس ۱۸۰ درایه مرتبط با ۱۸۰ مقدار فروش دارد. برای پیش بینی فروش هر ماه از پنج متغیر زیر استفاده می شود:

۱. فروش یک ماه قبل؛ ۲. فروش دو ماه قبل؛ ۳. فروش سه ماه قبل؛ ۴. میانگین متحرک^۱ سه ماهه و ۵. میانگین متحرک شش ماهه.

ماتریس دوم به این صورت ایجاد می شود که در ستون اول مقادیر فروش را درج نموده و در ستون های ۲ تا ۶ پنج متغیر بالا آورده می شود. از آنجا که ۱۸۰ داده فروش وجود داشته و متناظر هر داده پنج متغیر پیش بینی کننده موجود است، ماتریس به دست آمده، یک ماتریس ۶ در ۱۸۰ خواهد بود.

پس از حذف سطرهایی که دارای درایه ای خالی هستند، اندازه ماتریس حاصل شده ۶ در ۱۵۰ خواهد بود.

کل داده های موجود را مطابق شکل ۴ تقسیم نموده و تعریف هر یک از مجموعه ها و دلایل تقسیم بندی مزبور در ادامه آورده شده است.



شکل ۴. مجموعه داده ها

مجموعه داده تعلیم:

مجموعه داده تعلیم داده‌هایی هستند که در برآزش مدل به کار می‌روند. در این تحقیق حدوداً ۶۰ درصد داده‌های ابتدایی به عنوان مجموعه داده تعلیم انتخاب شده است. مجموعه داده اعتبارسنجی:

باید در انتخاب مدل، مقادیری برای پارامترهای مدل تعیین شود؛ برای نمونه، در الگوریتم SVM همان‌طور که در قبل توضیح داده شد، باید مقادیر بهینه‌ای برای سه پارامتر C ، Σ و γ انتخاب شود زیرا تغییر این پارامترها تأثیر زیادی بر نتایج کار و دقت پیش‌بینی می‌گذارد. روشن است اگر تنها از مجموعه داده‌های تعلیم برای تعیین مقادیر بهینه C ، Σ و γ استفاده شود، این خطر وجود دارد که مقادیر انتخاب‌شده تنها مقادیر بهینه برای همان مجموعه باشند؛ به سخن دیگر، با این کار خطا در نمونه پایین آمده، اما خطای عمومی کاهش نیافته است. به این ترتیب، برای انتخاب پارامترهای بهینه، مدل نیازمند داده‌هایی است که در برآزش مدل از آنها استفاده شده است. به این داده‌ها مجموعه داده اعتبارسنجی گفته می‌شود که در تحقیق حاضر حدوداً ۲۰ درصد از کل داده‌ها را به عنوان مجموعه داده اعتبارسنجی انتخاب شده است.

به این روش انتخاب پارامترها، اعتبارسنجی ضربدیری^۱ گفته می‌شود و گام‌های اجرای آن بعد از تفکیک داده‌ها به صورت زیر است:

۱. مجموعه‌هایی از مقادیر پارامترها را انتخاب می‌کنیم؛ برای نمونه، در مثال SVM ۱۰ مجموعه از مقادیر C ، Σ و γ را به صورت زیر انتخاب می‌شود:

$$P1: (C_1, \Sigma_1, \gamma_1)$$

$$P2: (C_2, \Sigma_2, \gamma_2)$$

$$P3: (C_3, \Sigma_3, \gamma_3)$$

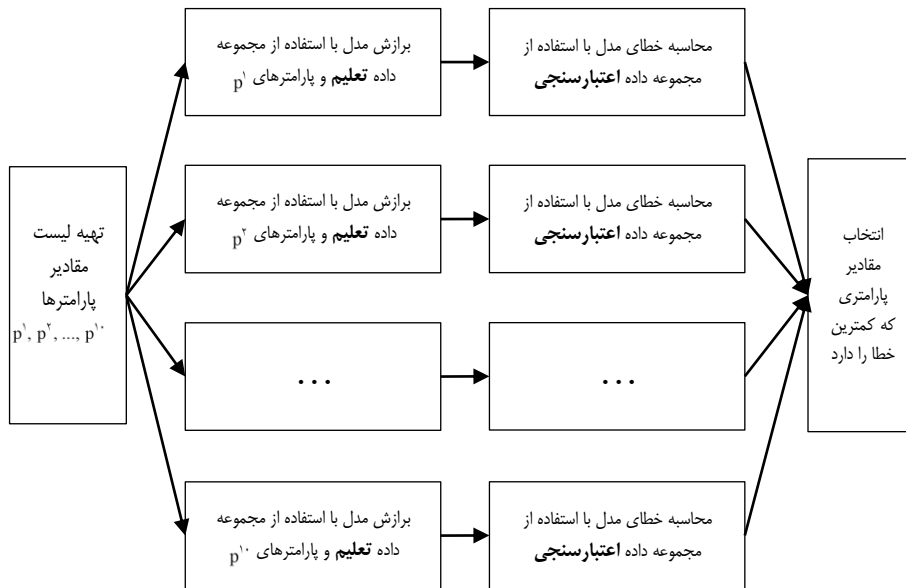
$$P10: (C_{10}, \Sigma_{10}, \gamma_{10})$$

۲. با استفاده از مجموعه داده‌های تعلیم، مدل در هریک از این ۱۰ حالت برآزش می‌شود، در نتیجه به ۱۰ مدل متفاوت دست یافته می‌شود.

۳. خطای این ۱۰ مدل را با مجموعه داده‌های اعتبارسنجی مورد سنجش قرار داده می‌شود.

۴. مجموعه مقادیر پارامتری را که کمترین خطا را داشته است، به عنوان پارامترهای بهینه مدل انتخاب می‌گردد.

در شکل ۵ الگوریتم اجرای گام‌های فوق مشاهده می‌شود:



شکل ۵. مراحل اجرایی

در عمل ممکن است برای هریک از پارامترها ۱۰ مقدار متفاوت انتخاب شده و تمامی مجموعه‌های سه‌تایی از آن‌ها به عنوان P های شکل بالا در نظر گرفته شود؛ در این صورت در کل $10^3 = 1000$ مجموعه خواهیم داشت و انتخاب پارامتر بهینه از میان این مجموعه‌ها انجام می‌گیرد.

مجموعه داده آزمون. در نهایت برای ارزیابی عملکرد مدل، به مجموعه داده‌های جدیدی نیاز می‌باشد. با توجه به اینکه مدل پیش‌بینی مطلوب باید عملکرد خوبی روی داده‌های جدید داشته باشد. بنابراین اگر از مجموعه داده‌های تعلیم یا اعتبارسنجی برای سنجش عملکرد مدل تحقیق استفاده شود، سنجش درست نبوده و تخمین ما از عملکرد مدل بهتر از عملکرد واقعی آن خواهد بود، زیرا در برازش مدل از دو مجموعه ذکر شده استفاده شده و به اصطلاح، مدل داده‌های این دو مجموعه را دیده است.

به این ترتیب، به مجموعه داده‌هایی نیاز داریم که در هیچ مرحله‌ای از آن استفاده نشده باشد. به این مجموعه، مجموعه داده آزمون گفته می‌شود و تنها یک بار در انتهای کار و برای

سنجش عملکرد واقعی مدل (سنجش خطای عمومی^۱) از آن استفاده می‌شود. در این مقاله حدوداً ۲۰ درصد داده‌های انتهایی را به‌عنوان مجموعه داده تعلیم انتخاب گردیده است.

۴. تجزیه و تحلیل داده‌ها

برای ارزیابی عملکرد مدل پیشنهادی این تحقیق، دقت پیش‌بینی مدل را با مدل‌های رگرسیون خطی و مدل SVR ساده مقایسه شده است. لازم به ذکر است که قبل از اجرای رگرسیون خطی، برای از بین بردن همبستگی متغیرها از آنالیز اجزای اصلی استفاده گردید و پارامترهای بهینه مدل SVR ساده را از طریق اعتبارسنجی ضربدری به دست آمده است. همچنین برای محاسبه خطای مدل‌ها از مجذور میانگین مربعات خطا^۲ استفاده شده است. نحوه محاسبه این خطا را در رابطه ۹ مشاهده می‌شود:

$$RMSE = \sqrt{\frac{\sum_{i=1}^n (T_i - P_i)^2}{n}} \quad \text{رابطه ۹}$$

که در آن T مقدار حقیقی، P مقدار پیش‌بینی شده و n تعداد کل نقاط می‌باشد. عملکرد این سه مدل در جدول ۱ مشاهده می‌شود. نتایج نشان می‌دهد که مدل پیشنهادی این تحقیق کمترین خطای عمومی را از میان مدل‌های دیگر دارا بوده و بیانگر عملکرد بهتر این مدل می‌باشد.

نتایج		
	خطا در نمونه	خطای عمومی
رگرسیون خطی	۰/۳۶۱۰	۰/۸۴۰۶
SVR	۰/۴۵۰۹	۰/۷۵۷۰
مدل پیشنهادی	۰/۵۷۳۷	۰/۶۳۹۹

جدول ۱. مقایسه عملکرد مدل‌ها

1. Generalization error
1. Root mean square error (RMSE)

۵. نتیجه‌گیری

در مقاله حاضر، پیش بینی فروش کاشی و سرامیک به کمک رگرسیون بردار پشتیبان انجام شده و نشان داده شده است که پیش بینی بوسیله رگرسیون مزبور در مسائل کاربردی از نتایج دقیق تری برخوردار می باشد. از آنجایی که وجود اغتشاش در داده ها نتایج پیش بینی را به مقدار زیادی تحت تأثیر قرار می دهد، لذا پیش پردازش داده ها از اهمیت زیادی برخوردار است. در این پژوهش با کاهش بعد داده ها و با پیشنهاد استفاده از آنالیز مولفه های مستقل و یادگیری منیفلد تا حد قابل ملاحظه ای اغتشاشات موجود در داده های مورد مطالعه کاهش یافته است.

از جمله روش های کاهش خطای عمومی در مقابل کاهش خطا در نمونه، استفاده از روش اعتبار سنجی ضربدری است که در مقاله حاضر از این روش جهت دستیابی به پارامترهای بهینه مدل استفاده شده است. همانگونه که در جدول شماره ۱ نشان داده شده است، مدل پیشنهادی این پژوهش برای پیش بینی در مقایسه با برخی از مدل های بررسی شده، دلیل دارا بودن خطای عمومی کمتر از دقت بالاتری برخوردار می باشد. بدین ترتیب برای پیش بینی تقاضا در صنعت کاشی و سرامیک توصیه می شود از مدل پیشنهادی این پژوهش استفاده شود.

پیشنهاد برای تحقیقات آتی. برای تحقیقات آتی موارد زیر پیشنهادی می شود:

- بررسی مدل با اضافه کردن متغیرهای دیگر به متغیرهای پیش بینی مدل تحقیق مثل اطلاعات رقبا، اطلاعات محیط صنعت؛

- بهره گیری از مفاهیم آمار بیزی در طراحی مدل های پیش بینی؛

- طراحی مدل های جدید با مجموعه ای از داده های بیشتر برای زمان های دورتر و برای سایر صنایع.

منابع

1. Delaney, R., & R. Wilson. (2003). *14th Annual State of Logistics Report*.
2. Simchi-Levi, D. (2004). *Managing The Supply Chain*. New York :McGraw-Hill
3. Henkoff, R. (1994). Delivering the Goods. *Fortune*, 64-78
4. Tay, F.E.H., & Cao, L.J. (2003). Support vector machine with adaptive parameters in financial time series forecasting. *IEEE Transactions on Neural Networks* 14, 1506-1518.
5. Tay, F.E.H., & Cao, L.J. (2001). Improved financial time series forecasting by combining support vector machines with self-organizing feature map. *Intelligent Data Analysis* 5, 339-354.
6. Chang, P.C., Wang, Y.W., & Tsai, C.Y. (2005). Evolving neural network for printed circuit board sales. *Expert Systems with Applications* 29 (1), 83-92.
7. Vapnik, V.N. (1999). An overview of statistical learning theory. *IEEE Transactions on Neural Networks* 10, 988-999.
8. Vapnik, V.N. (2000). *The Nature of Statistical Learning Theory*. New York :Springer.
9. Chang, P.C., Wang, Y.W., & Liu, C.H. (2006). Combining SOM and GA-CBR for flow time prediction in semiconductor manufacturing factory. *Lecture Notes in Computer Science* 4259, 767-775.
10. Wang, Y.W., Liu, C.H., & Fan, C.Y. (2009). The hybrid model development of clustering and back propagation network in printed circuit board sales forecasting. *Opportunities and Challenges for Next-Generation Applied Intelligence* 214, 213-218.
11. Chi, J.L., & Yen, W.W. (2010). Combining independent component analysis and growing hierarchical self-organizing maps with support vector regression in product demand forecasting. *Int. J. Production Economics* 128, 603-613.
12. Back, A., & Weigend, (1997). Discovering structure in finance using independent component analysis. *Proceedings of the Fifth International Conference on Neural Networks in Capital Market*, Kluwer Academic, pp. 15-17.
13. Kiviluoto, K., & Oja, E. (1998). Independent component analysis for parallel financial time series. *Proceedings of the Fifth International Conference on Neural Information*, Tokyo, Japan, pp. 895-898.
14. Oja, E., Kiviluoto, K., & Malaroiu, S. (2000). Independent component analysis for financial time series. *Proceeding of the IEEE 2000 Adaptive Systems for Signal Processing, Communications, and Control Symposium*, Lake Louise, Canada, pp. 111-116.
15. Cao, L.J. (2003). Support vector machines experts for time series forecasting. *Neurocomputing* 51, 321-339.

16. Lai, R.K., Fan, C.Y., Huang, W.H., & Chang, P.C. (2009). Evolving and clustering fuzzy decision tree for financial time series data forecasting. *Expert Systems with Applications* 36 (2 PART 2), 3761–3773.
17. Hung, C., & Tsai, C.F. (2008). Market segmentation based on hierarchical self-organizing map for markets of multimedia on demand. *Expert Systems with Applications* 34, 780–787.
18. Mohandes, M.A., Halawani, T.O., Rehmam, S., & Hussain, A.A (2004). Support vector machines for wind speed prediction. *Renewable Energy* 29, 939–947.
19. Thissen, U., Brakel, R.V., De Weijer, A.P., Melssen, W.J., & Buydens, L.M.C. (2003). Using support vector machines for time series prediction. *Chemometrics and Intelligent Laboratory Systems* 69, 35–49.
20. Lu, C.-J., Lee, T.-S., & Chiu, C.-C. (2009). Financial time series forecasting using independent component analysis and support vector regression. *Decision Support Systems* 47 (2), 115–125.
21. Pai, P.F., Yang, S.L., & Chang, P.T. (2009). Forecasting output of integrated circuit industry by support vector regression models with marriage honey-bees optimization algorithms. *Expert Systems with Applications* 36 (7), 10746–10751.
22. Yang, Y., Fuli, R., Huiyou, C., & Zhijiao, X. (2007). SVR mathematical model and methods for sale prediction. *Journal of Systems Engineering and Electronics* 18 (4), 769–773.
23. He, W., Wang, Z., & Jiang, H. (2008). Model optimizing and feature selecting for support vector regression in time series forecasting. *Neurocomputing* 72 (1–3), 600–611.
24. Fang, S.F., Wang, M.P., Qi, W.H., & Zheng, F. (2008). Hybrid genetic algorithms and support vector regression in forecasting atmospheric corrosion of metallic materials. *Computational Materials Science* 44 (2), 647–655.