



فصلنامه علمی پژوهشی دانش سرمایه‌گذاری
سال دهم / شماره چهارم / زمستان ۱۴۰۰

پیش‌بینی روند روزانه قیمت سهام با استفاده از متن کاوی احساسات کاربران شبکه اجتماعی و داده‌کاوی نماگرهای تکنیکال

کامل ابراهیمیان

دانشجوی دکتری مدیریت مالی، گروه مدیریت مالی، واحد قزوین، دانشگاه آزاد اسلامی، قزوین، ایران
Kamel61tb@gmail.com

ابراهیم عباسی

دانشیار و عضو هیئت علمی دانشگاه الزهراء، دانشکده علوم اجتماعی و اقتصادی، ده ونک، تهران، ایران (نویسنده مسئول)
AbbasiEbrahim2000@alzahra.ac.ir

اکبر عالم تبریز

استاد، گروه مدیریت صنعتی، دانشکده مدیریت و حسابداری، دانشگاه شهید بهشتی، تهران، ایران
Tabriz-a@sbu.ac.ir

امیر محمدزاده

دانشیار، گروه مدیریت مالی، واحد قزوین، دانشگاه آزاد اسلامی، قزوین، ایران
A.Mohammadzadeh@qiau.ac.ir

تاریخ دریافت: ۱۴۰۰/۰۱/۲۸ تاریخ پذیرش: ۱۴۰۰/۰۲/۰۴

چکیده

این پژوهش به پیش‌بینی روند آتی قیمت سهام در بازه کوتاه مدت روزانه با استفاده از تحلیل نظرات سرمایه‌گذاران در شبکه‌ی اجتماعی سهامیاب می‌پردازد. قابلیت پیش‌بینی پذیری بازارهای سهام، به خاطر دارا بودن یک سیستم پیچیده، پویا و غیرخطی همواره از چالش‌های پژوهشگران بوده است در این تحقیق، برای اولین بار، با تحلیل احساسات نظرات کاربران و ترکیب آن با ۲۰ نماگر تکنیکال به کمک سه الگوریتم درخت تصمیم، بیز ساده و ماشین بردار پشتیبان، مدلی با دقت ۷۲/۰۸ درصد برای پیش‌بینی جهت حرکت سهام توسعه‌یافت و به پیش‌بینی روند کوتاه مدت سهام پرداخته شد بر اساس نتایج، ماشین بردار پشتیبان، عملکردی بهتر از دو الگوریتم دیگر از خود نشان داد. همچنین مشخص شد حجم معاملات روز آتی و تعداد نظر ها دارای همبستگی معناداری است و نتایج آزمون علیت گرنجر نشان داد می‌توان برای پیش‌بینی قیمت سهام، از تجمیع احساسات روزانه کاربران نیز بهره جست.

واژه‌های کلیدی: پیش‌بینی قیمت سهام، الگوریتم‌های طبقه‌بندی، تحلیل احساسات، علیت گرنجر.

۱- مقدمه

در گذشته همیشه این سؤال مطرح بوده که تا چه میزان می‌توان از قیمت و رفتار گذشته سهام برای پیش‌بینی آن در آینده استفاده کرد؟ (فاما ۱۹۶۵). هدف این پژوهش، پیش‌بینی روند حرکت سهام با استفاده از تحلیل احساسات و نماگرهای تکنیکال می‌باشد بررسی همبستگی و رابطه‌ی بین حجم نظرات با حجم معاملات روز بعد و بررسی علت گرنجری تجمیع احساسات روزانه بر نسبت تغییر قیمت روز بعدی و نیز تحلیل احساسات کاربران وب سایت سهام یاب از دیگر اهداف اصلی پژوهش می‌باشد بر همین اساس از سه الگوریتم درخت تصمیم^۱ (DTREE)، ماشین بردار پشتیبان^۲ (SVM) و بیز ساده^۳ (NBAYES) در این پژوهش استفاده می‌شود. نتایج حاصل نشان دهنده دقت قابل قبول برای پیش‌بینی جهت حرکت سهام است و نوآوری این پژوهش در پیش‌بینی روند کوتاه مدت سهام با ترکیب تحلیل حساسیت نظر کاربران و استفاده از ۲۰ نماگر شاخص‌های علم تکنیکال است.

در بخش دوم مبانی نظری تحقیق، روش‌ها و الگوریتم‌های به کار رفته در پژوهش توضیح داده می‌شوند. در بخش سوم روش شناسی پژوهش بیان شده و در بخش چهارم و پنجم به ترتیب نتایج حاصل از مدل سازی و بحث و نتیجه گیری بیان خواهد شد.

۲- مبانی نظری و پیشینه‌ی تحقیق

۱-۱- داده‌کاوی

روش‌های داده‌کاوی را می‌توان به دو دسته‌ی روش‌های باناظر و بدون ناظر تقسیم کرد. روش‌های باناظر در حقیقت شامل الگوریتم‌هایی است که وظیفه‌ی پیش‌بینی را بر عهده دارند. این الگوریتم‌ها با یادگیری رابطه‌ی بین متغیرهای ورودی و متغیرهدف، الگوی پنهان بین آن‌ها را شناسایی کرده و از این طریق پیش‌بینی را انجام می‌دهند. الگوریتم‌های طبقه بندی و رگرسیون از روش‌های باناظر داده‌کاوی محسوب می‌شوند. در روش‌های بدون ناظر متغیرهدف یا متغیرقابل پیش‌بینی وجود ندارد. خوشه بندی و قوانین انجمنی از روش‌های بدون ناظر داده‌کاوی محسوب می‌شوند در روش‌های باناظر داده‌ها به دو مجموعه آموزشی (۷۰ درصد) و آزمایشی (۳۰ درصد) افراز می‌شوند که با مجموعه‌ی آموزشی، مدل آموزش داده شده و با مجموعه‌ی آزمایشی صحت و دقت مدل را می‌سنجد.

۲-۲- متن‌کاوی

هدف از متن‌کاوی پیدا کردن دانش‌های نهفته شده در متون می‌باشد. فرایند متن‌کاوی دارای چهار مرحله اصلی می‌باشد. در ابتدا متون و نظرات افراد جمع‌آوری می‌شود. در مرحله دوم متون پاک سازی و عملیات پیش پردازش متون بصورت زیر انجام خواهد شد.

۱) تبدیل کلمات محاوره ای به رسمی

۲) حذف کاراکترهای اضافه مثل علامت سؤال، علامت تعجب، کاراکترها و کلمات غیرفارسی

۳) یکسان سازی حروفی که دارای چندین شکل هستند مثل حرف ی، ک، ا، آ

۴) حذف کلمات stop word

۵) ریشه‌یابی

مرحله سوم تشکیل ماتریسی به عنوان Term-Document -Matrix است که ستون‌های این ماتریس کلمات یا اصطلاح‌های (Term) استخراج شده هستند و سطرهای این ماتریس نظرات یا سندها (Document) هستند. به عبارت دیگر این ماتریس نشان‌دهنده فراوانی هر کلمه در هر سند (متن، نظر و ...) است. برای به دست آوردن نتایج از الگوریتم وزن دهی به نام TF-IDF^۴ استفاده می‌شود که میزان اهمیت یک اصطلاح نسبت به یک سند، را نشان می‌دهد فرض کنید N تعداد کل اسناد، n_i تعداد اسنادی است که شامل اصطلاح i باشند، D_{ji} فراوانی اصطلاح i در سند j است. تعداد کل اصطلاحات سند j باشد. در این صورت اهمیت اصطلاح i در سند j از رابطه (۱) به دست می‌آید.

$$TF_i - IDF_j = \frac{D_{ji}}{D_j} \log \frac{N}{n_i} \quad (1)$$

در مرحله چهارم، برای تعیین قطبیت نظرات در دو طبقه مثبت و منفی، بخشی از نظرات را برچسب زده که نقش متغیر هدف را ایفا می‌کنند و سپس با اجرای گام‌های اول تا سوم، اصطلاحات، استخراج خواهند شد این اصطلاحات نقش متغیرهای ورودی را بازی خواهند کرد.

۲-۳- الگوریتم‌های داده‌کاوی

در این پژوهش از الگوریتم درخت تصمیم، بیز ساده و ماشین بردار پشتیبان (SVM) استفاده شده است

۲-۳-۱- الگوریتم درخت تصمیم

نحوه طبقه‌بندی درخت تصمیم ب، ب صورت است که به محض ورود یک نمونه ناشناخته (نمون، ی که طبقه یا برچسب آن مشخص نیست)، یک سری سؤالات به ترتیب در مورد آن نمونه پرسیده خواهد شد و با پاسخ گویی به هر سؤال مسیر حرکت در هر درخت مشخص می‌شود.

۲-۳-۲- الگوریتم بیز ساده

بیز ساده را می‌توان یک مدل بر مبنای احتمال شرطی^۵ در نظر گرفت. فرض کنید $X = (x_1, x_2, \dots, x_n)$ برداری از n ویژگی را بیان کند، که به صورت متغیرهای مستقل عمل می‌کنند. ب، ب ترتیب می‌توان احتمال رخداد C_k یعنی $P(C_k | x_1, x_2, \dots, x_n)$ به عنوان یکی از حالت‌های کلاس رخدادها، مختلف به ازای k های متفاوت از رابطه (۲) که همان قضیه بیز است نمایش داد.

$$P(C_k | X) = P\left(C_k \mid x_1, x_2, \dots, x_n\right) = \frac{P(C_k)P(X|C_k)}{P(X)} \quad (2)$$

۲-۳-۳- الگوریتم ماشین بردار پشتیبان (SVM)

الگوریتم ماشین بردار پشتیبان (SVM) توسط واپینگ و همکاران^(۱۹۹۵) معرفی شد بر این فرض استوار است که هدف مسأله، طبقه بندی بین دو دسته است. SVM را می‌توان پدید آورنده‌ی خط یا ابرصفحه در بین مجموعه داده‌ها و به منظور طبقه بندی آن‌ها در نظر گرفت.

۲-۴- ارزیابی مدل

کارایی مدل با استفاده از ماتریس درهم ریختگی مطابق جدول (۱) سنجیده می‌شود. این ماتریس چگونگی عملکرد الگوریتم دسته بندی را با توجه به مجموعه داده‌ی ورودی و نحوه‌ی تفکیک آن در دسته‌های مناسب و نامناسب مورد تحلیل قرار داده و نمایش می‌دهد. در جدول ۱، TP نمونه‌های مثبتی است که به درستی، مثبت تشخیص داده شده‌اند. TN نمونه‌های منفی که به درستی، منفی تشخیص داده شده‌اند. FN نمونه‌های مثبتی که به صورت نادرست، منفی تشخیص داده شده‌اند و FP نمونه‌های منفی که به صورت نادرست، مثبت تشخیص داده شده‌اند.

جدول (۱): ماتریس درهم ریختگی

		پیش‌بینی‌ها	
		درست T	نادرست F
نمونه‌های واقعی	مثبت P	TP	FN
	منفی N	FP	TN

برای سنجش دقت مدل نیز از سنجه‌ی نسبت کل پیش‌بینی‌های درست استفاده می‌کنیم که نحوه محاسبه آن در جدول (۲) بیان شده است.

جدول (۲): سنجش دقت مدل

نحوه محاسبه	شرح	نسبت
$TPR = \frac{TP}{TP + FN}$	نسبت نمونه‌هایی که به درستی مثبت تشخیص داده شده‌اند به تعداد کل نمونه‌های مثبت موجود	TPR نرخ مثبت درست
$TNR = \frac{TN}{FP + TN}$	نسبت نمونه‌هایی که به درستی منفی تشخیص داده شده‌اند به تعداد کل نمونه‌های منفی موجود	TNR نرخ منفی درست
$Accuracy = \frac{TP + TN}{TP + TN + FN + FP}$	نسبت کل پیش‌بینی‌های درست	Accuracy دقت کلی الگوریتم

۲-۵- همبستگی

تحلیل همبستگی برای مشخص کردن نوع و شدت رابطه‌ی دو متغیر کمی از رابطه (۳) محاسبه می‌شود:

$$\rho = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}} \quad (3)$$

مفهوم معنی داری همبستگی این است که همبستگی به دست آمده بین دو متغیر را می‌توان شانس و تصادفی در نظر گرفت یا این که واقعاً بین دو متغیر همبستگی وجود دارد آزمون فرضیه بصورت زیر:

$$\begin{cases} H_0: \rho = 0 & \text{همبستگی وجود ندارد} \\ H_1: \rho \neq 0 & \text{همبستگی وجود دارد} \end{cases}$$

و در سطح ۹۵ درصد سنجیده می‌شود اگر مقدار p-value کوچکتر از ۰/۰۵ باشد فرض صفر رد می‌شود یعنی همبستگی معنی داری بین دو متغیر وجود خواهد داشت.

۲-۶- علیت گرنجری^۷

گرنجر (۱۹۶۹) با استفاده از این واقعیت که آینده نمی‌تواند علت حال یا گذشته باشد بیان می‌کند اگر مقادیر جاری سری زمانی Y_t با استفاده از مقادیر گذشته سری X_t با دقت بیشتری نسبت به حالتی که از اطلاعات X_t استفاده نمی‌شود، پیش‌بینی شود، در این صورت X_t را علت گرنجری Y_t گویند برای آزمون فرضیه « X_t علت گرنجری Y_t نیست»، مدل اتورگرسیو برداری^۸ (VAR) از رابطه (۴)

$$Y_t = \sum_{i=1}^k \alpha_i Y_{t-i} + \sum_{i=1}^k \beta_i X_{t-i} + \mu_t \quad (4)$$

تعریف می‌شود. در این صورت اگر برای هر $i = 1, 2, \dots, k$ ، $\beta_i = 0$ باشد، در نتیجه X_t علت گرنجری Y_t نخواهد بود. جیوئیک^۹ (۱۹۸۴) بیان می‌کند که اعتبار آزمون به مانایی^{۱۰} هر دو سری زمانی بستگی دارد. یعنی هر دو سری زمانی می‌بایست مانا باشند. برای آزمون مانایی سری‌های زمانی می‌توان از آزمون دیکی فولر^{۱۱} استفاده نمود. این آزمون، فرضیه زیر را مورد بررسی قرار می‌دهد.

$$\begin{cases} H_0: \text{سری زمانی مانا نمی‌باشد} \\ H_1: \text{سری زمانی مانا می‌باشد} \end{cases}$$

۲-۷- پیشینه تحقیق

احساسات می‌تواند نوسانات کوتاه مدت بازار را به دنبال داشته باشند شوماکر^{۱۲} و چن^{۱۳} (۲۰۰۹) تأثیر اخبارهای فوری بر قیمت سهام را پس از بیست دقیقه از انتشار آن‌ها ارزیابی کردند. بولن و همکاران (۲۰۱۱) از داده‌های توئیت برای پیش‌بینی روند قیمت سهام استفاده کردند. آن‌ها با استفاده از الگوریتم GPOMS^{۱۴} به تحلیل احساسات ده میلیون توئیت پرداخته و از این طریق قیمت بسته شدن شاخص صنعتی داوجونز^{۱۵} را با دقت ۸۷/۶٪ پیش‌بینی کردند. میتال و گوئل^{۱۶} (۲۰۱۲) اساس کار خود را بر مطالعه‌ی بولن و همکاران (۲۰۱۱) قرار دادند اما با مجموعه‌ی بزرگتری از نظرات با اندازه‌ی ۴۰۰ میلیون و دقتی که آن‌ها به دست آوردند از ۸۷/۶٪ به ۷۵٪ کاهش یافت. لی و همکاران (۲۰۱۴) روشی را برای تعیین اهمیت متن‌کاوی^{۱۷} در پیش‌بینی سهام ارائه دادند. مشاری و همکاران (۱۳۹۹) در پژوهشی به بررسی قابلیت پیش‌بینی پذیری نقاط شروع (کف) و پایان (سقف) روند کوتاه مدت قیمت سهام با استفاده از مدل نایوبیز پرداختند. وطن پرست و همکاران (۱۳۹۸) از یک شبکه عصبی LM-BP براساس سری‌های زمانی قیمتی روشی برای پیش‌بینی قیمت سهام معرفی کردند.

۳- سوالات پژوهش

- آیا تغییرات جهت حرکت سهام با کمک تکنیک‌های داده‌کاوی و متن‌کاوی امکان‌پذیر است؟
- الگوریتم ارائه شده نسبت به سایر روش‌ها چه مزیت‌هایی دارد؟
- آیا قیمت‌های پیش‌بینی شده با دقت قابل قبولی پیش‌بینی پذیر هستند؟
- تاثیر شبکه‌های اجتماعی بر دقت پیش‌بینی نقاط ورود و خروج چقدر است؟

۴- روش شناسی و داده‌های پژوهش

هدف این پژوهش توسعه‌ی مدلی برای پیش‌بینی جهت حرکت سهام با رویکرد ترکیبی تحلیل احساسات و نماگرهای تکنیکالی است. بنابراین این پژوهش از نظر هدف، کاربردی و از نوع پژوهش‌های میدانی محسوب می‌شود.

۴-۱- داده‌های پژوهش

جامعه‌ی آماری اطلاعات جمع‌آوری شده از ۱۴ شرکت پذیرفته شده در بورس تهران می‌باشد. داده‌های متنی پژوهش نظرات و توئیت‌های روزانه‌ی کاربران وب سایت سهام‌یاب است. لیست کاملی از خواص جمع‌آوری شده اولیه را می‌توان در جدول (۲) مشاهده نمود. مدل سازی بر اساس ۲۰ نماگر تکنیکال قیمت‌های روزانه (جدول ۳) و یک متغیر تحلیل احساسات صورت گرفته است از داده‌های سال‌های ۱۳۹۷ تا ۱۳۹۸ جهت مدل سازی استفاده شده است.

جدول (۲) لیستی از مشخصه‌ها برای مجموعه

ماهیت	مشخصه‌ها برای مجموعه
نظرات	نام، نام کاربری، محتوای نظرات، تاریخ انتشار، زمان انتشار، منبع پیام
کاربران	نام، نام کاربری، نظرات کاربر،

جدول (۳) : نماگر های تکنیکال

شرح متغیر	عبارت انگلیسی	منبع
نوسان ساز مبتنی بر حرکت	Know Sure Thing	پرینگ، ۱۹۹۱
میانگین متحرک تطبیقی کافمن	Kaufman's Adaptive Moving Average	کافمن، ۲۰۱۳
منحنی کوپک	Coppock Curve	کوپک، ۱۹۶۲
شاخص ویلیام	Williams %R	ویلیامز، ۲۰۱۴
شاخص ورتکس	Vortex Indicator	بوتس و داگلاس، ۲۰۱۰
شاخص نوسان نگر لحظه، ی قیمت	Decision Point Price Momentum	مورفی، ۲۰۰۹
شاخص نوسان ساز نهایی	Ultimate Oscillator	ویلیامز، ۱۹۸۵
شاخص میانگین جهت دار	Average Directional Index	گوریب ایخلاص، ۲۰۱۸
شاخص گردش پول	Money Flow Index	دنگ و ساکارای، ۲۰۱۳
شاخص کانال کالا	Commodity Channel Index	باکر و همکاران، ۱۹۸۳
شاخص قدرت واقعی	True strength index	بالو، ۱۹۹۱
شاخص قدرت نسبی	relative strength index	مورفی، ۲۰۰۹
شاخص حجم منفی	Negative Volume Index	دیزارت، ۱۹۸۳
شاخص جریان پول چالکین	Chalkin Money Flow	نیمان، ۲۰۰۹
شاخص توده	MASS Index	دروسی، ۱۹۹۲
تغییر قیمت	Price Rate of Change	اسماعیل حسین، ۲۰۱۸
برد واقعی میانگین	Average True Range	ویلدنر، ۱۹۷۸
اسیلاتور نوسان قیمت	percentage price oscillator	انتاکاریس، ۲۰۱۹
استوکاستیک RSI	STOCHASTIC RSI	چانده، ۱۹۹۴
استوکاستیک %D	Stochastic oscillator	پرسون، ۲۰۰۴

۲-۴- روش انجام پژوهش

۲-۴-۱- جمع‌آوری و پیش‌پردازش داده‌ها

نماد شرکت‌های استفاده شده همراه با تعداد نظرات جمع‌آوری شده در بازه‌ی زمانی ۱۳۹۷/۱/۱ الی ۱۳۹۸/۱۲/۲۹ از وب سایت سهام یاب در جدول (۲) آورده شده است. تعداد کل توییتها این ۱۴ سهم ۵۵۰،۶۷۴

نظر بوده و با حذف نظرات تکراری به ۵۱۴،۷۰۵ نظر کاهش یافته است، برای پیش‌پردازش نظرات کاربران، مراحل زیر لحاظ شده است:

(۱) نظرات تکراری و نظرات آخر هفته‌ها و تعطیلات رسمی و زمان توقف نماد، حذف شدند.

(۲) در نظرات علامت @ را با متن "اتساین" جایگزین شد.

(۳) #نماد هر سهم در نظرات با کلمه "نمادسهم" جایگزین شد.

برای به دست آوردن نماگرها از نرم افزار TSE Client، آمی پروکر و اکسل و از طریق برنامه نویسی استفاده شده است و به از زبان برنامه نویسی R برای داده‌کاوی و متن‌کاوی استفاده شده است.

۴-۲-۲- برچسب زنی نظرات کاربران

در این پژوهش به جای استفاده از واژه‌نامه‌های احساسی از تحلیل قطبیتی احساسات سرمایه‌گذاران استفاده می‌شود چون اگرچه در تعیین احساسات توییت، ذهنیت و عینیت در نظر گرفته می‌شود، اما در بسیاری از موارد، توییت‌های بیطرف (عینی) در شرایط نرمال باید به عنوان توییت‌های قطبی در این پژوهش تلقی شوند. به عنوان مثال، این نظر "#خودرو هیچ خبری نیست فقط وقت تلف کردنه"، به احتمال زیاد به عنوان جمله‌ای خنثی و بی طرف علامت گذاری خواهد شد و بنابراین از تجزیه و تحلیل قطبیتی خارج می‌شود، در حالی که در واقعیت نشان می‌دهد یک کاربر انتظار دارد قیمت سهام کاهش یابد. ثانیاً در حوزه زبان فارسی برای نظرکاوی بازارهای مالی و حتی سایر علوم فارسی، واژه‌نامه احساسی و نرم افزارهای خاص وجود ندارد برای غلبه بر این مشکل، از برچسب زنی بخشی از توییت‌ها به کمک خبرگان برای آموزش الگوریتم‌های طبقه بندی استفاده نمودیم. در جدول (۴) برخی از توییت‌ها و برچسب‌های نوعی آنها و همچنین استدلال‌های اینکه چرا آنها ب، ین ترتیب برچسب گذاری شده‌اند بیان شده است.

جدول (۴) برخی نظریات کاربران و برچسب‌های آنها

برچسب	استدل	نظر
مثبت	نظری کوتاه همراه با احساسات مثبت نسبت به نماد #خودرو که روند رو به رشدی را طی خواهد کرد	#خودرو فردا تا ساعت ۱۱ هر چقدر حقوقی عرضه کرد بخرید که بعدش حالا حالاها نمی‌توانید بخرید.
خنثی	کاربر بدون اینکه احساسی از خود بروز دهد، در مورد تاثیر پروژه‌ها در شرایط حال و تفاوت آن نسبت به گذشته سوال می‌پرسد	#کیسون، در گذشته هم پروژه کم نداشته، پس چرا قبل از این هیچ تاثیری روی قیمت سهامش نگذاشته؟
منفی	کاربر نظر خود را در مورد رشد کاذب سهم اعلام می‌نماید، و تعویض سهم را پیشنهاد می‌کند.	#شپلی، این سهم ۱۰۰ تومانی با سفته بازی رفت بالا... سهم خود را تعویض کنید... از دواج # نکنید

۴-۲-۳- مدل سازی

در بخش اول پژوهش، مدلی برای طبقه‌بندی نظرات در دو طبقه‌ی مثبت و منفی توسعه داده می‌شود و پس از آموزش مدل و مشخص شدن کلاس تمامی نظرات، متغیرتجمیع احساسات روزانه محاسبه می‌شود. این متغیر از جمع جبری احساسات روزانه افراد به دست آمده است و عبارت است از تفاضل تعداد نظرات منفی از تعداد نظرات مثبت سهام مربوطه. به عبارت دیگر اگر نام یکی از سهام‌های موجود در جدول (۵) باشد در این صورت تجمیع احساسات روزانه این سهام از رابطه‌ی (۵) به دست می‌آید.

$$S_{ij} = S_{ij}^+ - S_{ij}^- \quad (5)$$

که در آن S_{ij}^+ تعداد نظرات مثبت سهام i در روز j ام، S_{ij}^- تعداد نظرات منفی همان سهام در روز j ام و S_{ij} تجمیع احساسات است. اگر $S_{ij} > 0$ باشد به این معنی است که کاربران نظر مساعدی نسبت به آن سهام دارند و احتمال می‌رود در روز بعدی جهت حرکت سهام رو به افزایش باشد و اگر $S_{ij} < 0$ باشد آن گاه عکس مطلب فوق می‌تواند برقرار باشد.

بخش دوم تحلیل مربوط می‌شود به رابطه همبستگی، علیت گرنجری و بخش سوم مربوط به مدل سازی با رویکرد ترکیبی براساس نماگرهای تکنیکالی و تحلیل احساسات برای پیش‌بینی جهت حرکت سهام است. بنابراین در این بخش ۲۱ متغیر ورودی روزانه برای یک سهم وجود دارد. ۲۰ متغیر مربوط به نماگرهای روزانه‌ی تکنیکالی و متغیر دیگر، تجمیع احساسات روزانه آن سهام است. متغیر هدف در این بخش جهت حرکت سهام است که در دو طبقه‌ی افزایشی و کاهش‌ی قرار می‌گیرد. اگر قیمت روز بعد سهام بیشتر از قیمت امروز باشد، در نتیجه امروز در طبقه افزایشی قرار می‌گیرد و بالعکس، اگر کمتر باشد، طبقه‌ی امروز کاهش‌ی خواهد بود.

آمار تعداد نظرات منتشر شده توسط کاربران، تعداد روزهای معاملاتی، رتبه دنبال کننده در سایت سهامیاب و نیز متوسط بازدید روزانه آن در جدول (۵) توصیف شده است.

جدول (۵): لیست سهام به همراه تعداد نظرات

نام سهام	تعداد نظرات	روز معاملاتی	رتبه دنبال کننده	متوسط بازدید روزانه	نام سهام	تعداد نظرات	روز معاملاتی	رتبه دنبال کننده	متوسط بازدید روزانه
اخابر	۴,۷۲۶	۴۳۵	۶۵	۲,۴۵۰	خودرو	۲۰,۸۴۴	۴۳۵	۶	۴۵,۵۸۱
حفاری	۱۰,۲۵۲	۳۴۹	۴۹	۳,۸۵۱	کیسون	۲۲,۴۸۴	۴۷۲	۱۱۷	۲,۷۳۰
فاراک	۱۱,۴۲۴	۴۳۶	۲۹	۴,۵۲۳	خساپا	۲۳,۷۹۷	۴۶۱	۵	۲۰,۶۲۴
شتران	۱۲,۱۳۹	۴۴۵	۱۳	۱۸,۱۶۶	وتجارت	۳۶,۰۹۱	۳۵۶	۲	۲۱,۸۵۵
ویصادر	۱۴,۳۳۷	۴۶۱	۱۲	۱۱,۲۹۵	ذوب	۳۸,۲۱۹	۴۶۳	۴	۷,۲۷۵
شبندر	۱۴,۳۶۸	۴۶۰	۱۰	۱۰,۴۶۱	تپکو	۴۲,۰۱۲	۴۵۸	۱۸	۳,۲۳۹
ستران	۱۷,۹۲۴	۴۵۳	۶۲	۲,۴۱۶	شپلی	۷۸,۰۱۵	۴۳۶	۲۳	۴,۷۷۸

۵- تجزیه و تحلیل

سه نوع تحلیل در این پژوهش مورد بررسی قرار گرفت اولین تحلیل مربوط به رویکرد پردازش زبان طبیعی و تحلیل احساسات، دومین تحلیل مربوط به همبستگی و علیت گرنجری و سومین تحلیل مربوط به رویکرد ترکیبی پیش‌بینی جهت حرکت سهام است.

۵-۱- پردازش زبان طبیعی و تحلیل احساسات

یک نمونه‌ی تصادفی به حجم ۲،۳۴۴ توئیت به عنوان نمونه‌های آموزشی و آزمایشی انتخاب گردید که به صورت دستی توسط خیرگان و کارشناسان بازار مالی، قطبیت آن‌ها در دو رده‌ی مثبت و منفی مشخص شد. از این تعداد ۱۶۴۰ توئیت (۷۰ درصد) به عنوان نمونه‌ی آموزشی و ۷۰۴ توئیت (۳۰ درصد) به عنوان نمونه‌ی آزمایشی انتخاب شدند. تعداد نظرات مثبت و منفی در نمونه‌ی آموزشی به ترتیب برابر با ۸۲۵ و ۸۱۵ و همین اعداد در نظرات آزمایشی برابر با ۳۷۵ و ۳۲۹ است. قبل از بیان نتایج الگوریتم‌های طبقه‌بندی تشریح می‌گردد که برای ساخت ویژگی‌ها (متغیرها) از unigram (تک کلمه ای‌ها) و bigram (دو کلمه ای‌ها) استفاده شده است. از الگوریتم‌های بیز ساده، درخت تصمیم و ماشین بردار پشتیبان برای ساخت مدل پیش‌بینی طبقه‌ی نظرات استفاده شده و نتایج ارزیابی آن‌ها برای نمونه‌های آموزشی و آزمایشی به ترتیب در جداول (۶) و (۷) آمده است. همان‌طور که از جدول (۶) مشخص است الگوریتم ماشین بردار پشتیبان در نمونه‌های آموزشی دارای بیشترین دقت با ۸۰/۳۷ درصد است و سپس الگوریتم بیز ساده دارای بیشترین دقت ۷۳/۶۶٪ بوده است. باید توجه شود که الگوریتم درخت تصمیم عملکرد بسیار ضعیفی از خود نشان داده است بنابراین نتایج آن به هیچ‌گونه عنوان قابل بررسی نمی‌باشد.

جدول (۶): نتایج ارزیابی الگوریتم‌های طبقه‌بندی نمونه‌های آموزشی (تحلیل احساسات)

الگوریتم	دقت کلی %	فراخوان توئیت‌های مثبت %	فراخوان توئیت‌های منفی %	ماتریس درهم ریختگی
SVM	۸۰/۳۷	۷۷/۳۳	۸۳/۴۴	منفی مثبت ۱۸۷ ۶۳۸ مثبت منفی ۶۸۰ ۱۳۵
DTREE	۶۲/۲۰	۳۵/۵۲	۸۹/۲۱	منفی مثبت ۵۳۲ ۲۹۳ مثبت منفی ۷۲۷ ۸۸
NBAYES	۷۳/۶۶	۶۹/۸۲	۷۷/۵۵	منفی مثبت ۲۴۹ ۵۷۶ مثبت منفی ۶۳۲ ۱۸۳

نتایج مشابهی نیز در ارزیابی نمونه‌های آزمایشی با این سه الگوریتم در جدول (۷) وجود دارد. با توجه به این جدول الگوریتم ماشین بردار پشتیبان دارای دقت ۷۰/۲۰ درصد است. لذا از این الگوریتم برای پیش‌بینی طبقه‌ی سایر نظرات استفاده خواهد شد.

جدول (۷): نتایج ارزیابی الگوریتم‌های طبقه‌بندی نمونه‌های آزمایشی (تحلیل احساسات)

الگوریتم	دقت کلی %	فراخوان توثیتهای مثبت %	فراخوان توثیتهای منفی %	ماتریس درهم ریختگی								
SVM	۷۰/۲۰	۷۲/۵۳	۶۷/۴۷	<table border="1"> <tr> <td>مثبت</td> <td>مثبت</td> <td>۲۷۲</td> <td>۱۰۳</td> </tr> <tr> <td>مثبت</td> <td>منفی</td> <td>۱۰۷</td> <td>۲۲۲</td> </tr> </table>	مثبت	مثبت	۲۷۲	۱۰۳	مثبت	منفی	۱۰۷	۲۲۲
مثبت	مثبت	۲۷۲	۱۰۳									
مثبت	منفی	۱۰۷	۲۲۲									
DTREE	۵۷/۳۹	۳۲/۲۷	۸۶/۰۲	<table border="1"> <tr> <td>مثبت</td> <td>مثبت</td> <td>۱۲۱</td> <td>۲۵۴</td> </tr> <tr> <td>منفی</td> <td>منفی</td> <td>۴۶</td> <td>۲۸۳</td> </tr> </table>	مثبت	مثبت	۱۲۱	۲۵۴	منفی	منفی	۴۶	۲۸۳
مثبت	مثبت	۱۲۱	۲۵۴									
منفی	منفی	۴۶	۲۸۳									
NBAYES	۶۵/۲۰	۶۰/۵۳	۷۰/۵۲	<table border="1"> <tr> <td>مثبت</td> <td>مثبت</td> <td>۲۲۷</td> <td>۱۴۸</td> </tr> <tr> <td>منفی</td> <td>منفی</td> <td>۹۷</td> <td>۲۳۲</td> </tr> </table>	مثبت	مثبت	۲۲۷	۱۴۸	منفی	منفی	۹۷	۲۳۲
مثبت	مثبت	۲۲۷	۱۴۸									
منفی	منفی	۹۷	۲۳۲									

پس از برچسب گذاری تمامی نظرات متغیرهای حجم روزانه‌ی نظرات، تعداد نظرات مثبت روزانه، تعداد نظرات منفی روزانه و تجمیع احساسات روزانه برای هر سهم محاسبه گردید.

۲-۵- همبستگی ها و علت گرنجری حجم توثیت ها

سؤال اساسی مهم این است که آیا حجم توثیتهای نماد x و حجم معاملات آن نماد در روز بعد رابطه و وجود دارد؟ و می‌توان انتظار داشت که بین حجم توثیتهای هر نماد با حجم معاملات روز بعد آن همبستگی معناداری وجود داشته باشد؟ مقدار ضریب همبستگی و p -value مربوط به آزمون فرض همبستگی در جدول (۸) بیانگر این است که در همه نمادها بجز خساپا و کیسون، بین حجم کل توثیت ها و حجم معاملات روز بعد، همبستگی مثبت و معناداری وجود دارد. اما همبستگی بین حجم توثیتهای مثبت و حجم معاملات روز بعد، فقط در ۹ نماد معنادار شده است و در آخر بین حجم توثیتهای منفی و حجم معاملات روز بعد، همبستگی معناداری مشاهده نمی‌شود.

با توجه به همبستگی بین حجم توثیت ها و حجم معاملات روز بعد، می‌توان یک گام فراتر رفته و این سؤال را مطرح نمود ک، یا تجمیع احساسات روزانه دارای اطلاعاتی برای پیش‌بینی نسبت تغییر قیمت سهام (PRC) می‌باشد؟ برای پاسخ ب، این سؤال می‌بایست از آزمون علیت گرنجر استفاده نمود اما قبل از آن، آزمون دیکی فولر برای بررسی مانایی هر دو سری زمانی الزامی می‌باشد. ذکر این نکته ضروری است که هر دو سری حداکثر باید در تأخیر یک مانا شوند. برای توضیح و بررسی‌های بیشتر، فرض کنید که علیت گرنجری تجمیع احساسات امروز (S_t) برای پیش‌بینی تغییر قیمت فردا (PRC_{t+1}) توسط مدل اتورگرسیو برداری (۶) تعریف شود.

$$PRC_{i+1} = \alpha + \alpha_1 PRC_i + \beta + \beta_1 S_i \quad (6)$$

بنابراین هر دو سری زمانی می‌بایست حداکثر در تأخیر یک مانا گردند و اگر نمادهایی نتوانند در هر یک از دو آزمون دیکی فولر و علیت گرنجر معنادار شوند، از تحلیل حذف خواهند شد. با توجه به جدول (۹) مشخص می‌شود که در تمامی نمادها آزمون دیکی فولر معنادار شده است و لذا هر دو متغیر PRC و S دارای شرایط مانایی در تأخیر یک هستند.

جدول (۸) همبستگی بین حجم توئیت ها و حجم معاملات روز بعد

نماد	کل توئیت ها		توئیت‌های مثبت		توئیت‌های منفی	
	P-value	همبستگی	P-value	همبستگی	P-value	همبستگی
اخابر	۰/۰۰۰۰	۰/۲۵۹۰	۰/۰۰۰۰	۰/۱۹۵۰	۰/۲۵۲۸	۰/۰۵۵۰
تپکو	۰/۰۰۰۱	۰/۲۱۴۰	۰/۰۲۳۳	۰/۱۲۲۰	۰/۱۷۰۸	۰/۰۷۳۶
تجارت	۰/۰۰۰۰	۰/۱۹۴۰	۰/۰۰۳۵	۰/۱۴۰۰	۰/۴۳۹۶	۰/۰۳۷۲
حفاری	۰/۰۰۰۰	۰/۲۰۰۰	۰/۰۰۰۰	۰/۱۹۳۰	۰/۸۲۱۳	-۰/۰۱۰۷
خسپا	۰/۰۵۴۱	۰/۰۸۹۹	۰/۵۶۰۲	۰/۰۲۷۲	۰/۱۷۸۱	۰/۰۶۲۹
خودرو	۰/۰۰۰۲	۰/۱۷۳۰	۰/۰۷۸۷	۰/۰۸۲۲	۰/۰۸۰۶	۰/۰۸۱۶
ذوب	۰/۰۳۰۴	۰/۱۰۲۰	۰/۴۹۱۸	۰/۰۳۲۴	۰/۱۰۵۰	۰/۰۷۶۴
ستران	۰/۰۰۰۶	۰/۱۶۵۰	۰/۱۰۴۹	۰/۰۷۸۰	۰/۱۰۱۱	۰/۰۷۸۸
شتران	۰/۰۰۰۰	۰/۲۰۹۰	۰/۰۰۰۱	۰/۱۷۶۰	۰/۹۶۱۰	۰/۰۲۲۶
شبندر	۰/۰۰۰۰	۰/۱۹۶۰	۰/۰۰۶۳	۰/۱۲۷۰	۰/۲۰۷۶	۰/۰۵۸۹
شپلی	۰/۰۰۰۳	۰/۱۹۲۰	۰/۰۰۴۲	۰/۱۵۲۰	۰/۸۳۱۶	۰/۰۱۱۳
فاراک	۰/۰۰۰۱	۰/۱۷۷۰	۰/۰۰۰۰	۰/۱۹۹۰	۰/۲۲۱۱	-۰/۰۵۷۰
کیسون	۰/۱۱۸۹	۰/۰۷۳۰	۰/۱۹۰۹	۰/۰۶۱۳	۰/۷۸۶۸	۰/۰۱۲۷
وبصادر	۰/۰۰۰۳	۰/۱۷۳۰	۰/۰۰۴۸	۰/۱۳۵۰	۰/۵۳۱۲	۰/۰۳۰۱

جدول (۹): نتایج آزمون دیکی فولر

نماد	نسبت تغییر قیمت		نماد	تجمیع احساسات		نماد	نسبت تغییر قیمت	
	آماره دیکی فولر	P-value		آماره دیکی فولر	P-value		آماره دیکی فولر	P-value
اخابر	-۱۲/۸۸	۰/۰۱*	ستران	-۱۴/۷۲	۰/۰۱*	اخابر	-۱۴/۴۶	۰/۰۱*
تپکو	-۱۲/۹۳	۰/۰۱*	شتران	-۱۱/۲۲	۰/۰۱*	تپکو	-۱۷/۱۸	۰/۰۱*
تجارت	-۱۴/۵۱	۰/۰۱*	شبندر	-۱۵/۱۹	۰/۰۱*	تجارت	-۱۵/۵۹	۰/۰۱*
حفاری	-۱۲/۵۰	۰/۰۱*	شپلی	-۱۵/۴۹	۰/۰۱*	حفاری	-۱۳/۰۸	۰/۰۱*
خسپا	-۱۲/۴۹	۰/۰۱*	فاراک	-۱۴/۵۲	۰/۰۱*	خسپا	-۱۷/۴۱	۰/۰۱*
خودرو	-۱۳/۰۶	۰/۰۱*	کیسون	-۱۵/۱۰	۰/۰۱*	خودرو	-۱۵/۰۱	۰/۰۱*
ذوب	-۱۴/۳۴	۰/۰۱*	وبصادر	-۱۵/۵۴	۰/۰۱*	ذوب	-۱۵/۰۷	۰/۰۱*

* کوچکتر از ۰/۰۱

با توجه به جدول (۱۰) آزمون علت گرنجری در تمامی نمادها بجز تپکو و تجارت معنادار شده است بنابراین این دو نماد از ادامه تحلیل حذف خواهند شد.

جدول (۱۰): نتایج آزمون علیت گرنجری تجمیع احساسات برای پیش‌بینی نسبت تغییر قیمت

نماد	آماره F	P-value	نماد	آماره F	P-value	نماد	آماره F	P-value	نماد	آماره F	P-value
اخابر	۲۷/۵۹	۰/۰۰۰۰	ستران	۳۴/۳۵	۰/۰۰۰۰	خساپا	۱۹/۰۲	۰/۰۰۰۰	فاراک	۶۴/۱۱	۰/۰۰۰۰
تپکو	۰/۸۴	۰/۳۵۸۷	شتران	۴۱/۶۰	۰/۰۰۰۰	خودرو	۵۰/۱۲	۰/۰۰۰۰	کیسون	۲۶/۶۵	۰/۰۰۰۰
تجارت	۰/۳۹	۰/۵۲۹۲	شبندر	۸/۱۱	۰/۰۰۴۶	ذوب	۲۱/۳۹	۰/۰۰۰۰	وبصادر	۷/۴۵	۰/۰۰۶۶
حفاری	۲۵/۵۴	۰/۰۰۰۰	شپلی	۳۳/۷۴	۰/۰۰۰۰						

۳-۵- رویکرد ترکیبی پیش‌بینی جهت حرکت سهام

بین احساسات کاربران و حجم معاملات و قیمت سهام در روز بعد رابطه و همبستگی وجود دارد. آمار توصیفی متغیرها (با حذف دو نماد تپکو و تجارت) و همچنین نقش آن‌ها در مدل سازی‌های این پژوهش در جدول (۱۱) نشان داده شده است.

جدول (۱۱): شاخص‌های آمار توصیفی متغیرهای مدل

نام متغیر	نقش	مینیمم	ماکسیمم	میانگین	انحراف معیار	چولگی	کشیدگی
VOLUME	ورودی	۲۰۰۰	۲۸۵۸۰۸۶۰۰۰	۶۳۶۰۸۱۷۰	۱۲۱۷۲۵۶۰۰	۷/۸۱۹۹	۹۹/۴۴۱۳
PRC	ورودی	-۱/۷۲۲۵	۰/۲۳۸۹	۰/۰۰۲۴	۰/۰۴۲۸	-۳۸/۱۶۶۸	۱۹۱۵/۲۱۴۳
ADX	ورودی	۳/۴۸۷۲	۹۶/۳۴۶۶	۲۹/۳۵۵۴	۱۷/۱۲۱۰	۰/۹۶۲۰	۰/۵۷۷۹
CCI	ورودی	-۶۴۷/۲۵۳۰	۴۵۹/۰۷۳۷	۲۹/۷۹۵۹	۱۱۳/۹۶۵۷	-۰/۱۸۹۲	۰/۵۶۵۴
CPPC	ورودی	-۱۲۲/۱۰۴۴	۱۹۱/۸۶۶۰	۱۰/۹۵۴۵	۳۱/۲۰۶۰	۱/۵۸۰۴	۹/۹۸۸۴
ATR	ورودی	۲/۵۷۱۴	۲۲۵۶	۱۹۹/۷۹۱۶	۲۳۰/۳۷۲۹	۳/۶۱۸۲	۲۱/۰۲۲۷
KAMMA	ورودی	۴۰۸/۱۰۸۷	۴۰۴۵۰/۵۴۰۰	۴۴۰۸/۳۲۶۰	۴۵۲۶/۴۳۸۰	۲/۹۰۵۳	۱۳/۶۶۸۸
KST	ورودی	-۵۹۳/۲۶۵۷	۱۴۹۳/۷۸۳۰	۹۶/۲۹۵۲	۲۱۶/۸۰۳۹	۲/۰۳۱۹	۹/۸۱۳۳
MASS	ورودی	۱۶/۹۹۶۵	۳۳/۷۱۷۴	۲۵/۳۵۶۹	۱/۵۹۲۶	-۰/۲۳۳۹	۳/۸۳۸۷
MFI	ورودی	۳/۷۳۷۹	۱۰۰	۵۹/۰۹۵۴	۱۹/۰۱۵۳	-۰/۲۶۷۲	-۰/۴۶۴۲
NVI	ورودی	۶۱۴/۳۷۰۹	۱۱۱۷/۵۶۰۰	۷۸۱/۳۸۰۳	۱۰۳/۰۰۵۸	-۰/۹۵۹۹	۱/۰۳۳۲
PPO	ورودی	-۲۶/۱۱۰۱	۲۳/۶۵۰۹	۱/۸۸۶۱	۵/۱۷۰۵	-۱/۴۷۴۶	۱۱/۱۷۷۲
RSI	ورودی	۱۳/۰۰۸۸	۹۸/۷۱۳۱	۵۵/۳۴۳۱	۱۴/۷۳۴۸	-۰/۱۲۶۰	-۰/۱۳۲۶
STOCHASTIC	ورودی	۰/۰۱۰۰	۱۰۰/۰۱۰۰	۵۴/۸۳۶۷	۳۰/۵۳۵۵	-۰/۱۶۷۲	-۱/۲۷۸۴
STOCHASTIC-۱۶	ورودی	۰/۰۰۰۰	۰/۹۹۹۹	۰/۵۰۲۴	۰/۳۷۱۶	-۰/۰۴۱۳	-۱/۵۳۱۱
TSI	ورودی	-۶۲/۷۶۸۹	۹۴/۴۵۵۲	۱۰/۶۶۰۲	۲۵/۹۶۷۸	۰/۰۱۵۳	۰/۰۸۶۰

نام متغیر	نقش	مینیمم	ماکسیمم	میانگین	انحراف معیار	چولگی	کشیدگی
UO	ورودی	۴/۸۱۴۴	۹۸/۹۰۰۳	۴۸/۷۰۱۳	۱۴/۵۷۵۷	۰/۵۲۶۶	۰/۶۰۷۴
VI	ورودی	۰/۵۴۸۹	۱/۵۲۱۸	۱/۰۱۳۰	۰/۱۵۴۲	۰/۰۵۲۹	-۰/۲۶۷۲
CMF	ورودی	-۰/۶۰۰۹	۰/۷۰۰۰	۰/۰۵۲۱-	۰/۲۰۸۶	۰/۱۹۴۷	-۰/۱۱۰۲
PMO	ورودی	-۱۷/۰۲۰۲	۳۳/۰۵۹۶	۳/۶۱۰۹	۶/۳۲۶۳	۰/۷۹۲۳	۳/۳۳۵۷
WILLIAMS SR	ورودی	-۱۰۰	۰	۴۵/۱۷۳۴-	۳۰/۵۳۵۵	-۰/۱۶۷۲	-۱/۲۷۸۴
MACD	ورودی	-۲۲۷۶/۱۳۱۱	۲۹۶۶/۱۰۰۰	۵۹/۸۵۹۳	۲۴۵/۰۷۱۲	۳/۳۳۰۲	۴۱/۴۴۸۳
S	ورودی	-۹۵	۱۰۵	۶/۳۷۸۸	۴۰/۷۴۹۵	-۰/۰۳۳۵	-۱/۱۷۹۷
(جهت حرکت سهام) Y	هدف	طبقه افزایشی: ۲۷۹۰			طبقه کاهش: ۲۵۴۵		

متغیرهدف، یک متغیررسته‌ای با دو دسته‌ی افزایشی و کاهش‌ی است و در کل اطلاعات این ۱۲ سهم ، ۵۳۳۵ روز را شامل می‌شوند که از این تعداد ۲۷۹۰ روز منجر به حرکتی افزایشی و ۲۵۴۵ روز نیز منجر به حرکت کاهش‌ی شده‌اند.

یک نمونه‌ی تصادفی به اندازه‌ی ۳۷۳۴ روز برای آموزش مدل‌های طبقه بندی در نظر گرفته شده است. نتایج آموزش و آزمایش مدل‌های طبقه بندی در جداول (۱۲) و (۱۳) گزارش شده است. با توجه به این دو جدول مشخص می‌شود که در نمونه‌های آموزشی و آزمایشی الگوریتم SVM دارای دقت بیشتری نسبت به دو الگوریتم دیگر می‌باشد و برخلاف مدل سازی تحلیل احساسات، الگوریتم بیز ساده عملکرد ضعیفی از خود نشان داده است. لذا از الگوریتم ماشین بردار پشتیبان برای رویکرد ترکیبی پیش‌بینی جهت حرکت سهام استفاده می‌شود و بر این اساس نتایج مدل سازی براساس این الگوریتم به صورت جداگانه برای هر یک از نمادها در جدول (۱۴) آورده شده است.

جدول (۱۲): نتایج ارزیابی الگوریتم‌های طبقه بندی نمونه‌های آموزشی (رویکرد ترکیبی)

الگوریتم	دقت کلی %	فراخوان مثبت %	فراخوان منفی %	خطا مثبت	خطا منفی	ماتریس درهم ریختگی
SVM	۷۵/۲۵	۷۶/۳۳	۷۴/۰۵	مثبت	۲۳/۶۷	مثبت
				منفی	۲۵/۹۵	منفی
DTREE	۶۹/۶۶	۶۹/۴۵	۶۴/۲۰	مثبت	۳۰/۵۵	مثبت
				منفی	۳۵/۸۰	منفی
NBAYES	۶۱/۴۶	۶۲/۳۷	۶۰/۴۴	مثبت	۳۷/۶۳	مثبت
				منفی	۳۹/۵۶	منفی

جدول (۱۳): نتایج ارزیابی الگوریتم‌های طبقه‌بندی نمونه‌های آزمایشی (رویکرد ترکیبی)

الگوریتم	دقت کلی %	فراخوان مثبت %	فراخوان منفی %	ماتریس درهم ریختگی		
				خطا	کاهش	افزایش
SVM	۷۲/۰۸	۶۸/۶۳	۷۵/۶۳	افزایش	۵۵۸	۳۱/۳۷
				کاهش	۱۹۲	۲۴/۳۷
DTREE	۶۴/۲۷	۵۸/۷۹	۶۹/۹۲	افزایش	۴۷۸	۴۱/۲۱
				کاهش	۲۳۷	۳۰/۰۸
NBAYES	۴۹/۴۱	۵۱/۹۱	۴۶/۸۳	افزایش	۴۲۲	۴۸/۰۹
				کاهش	۴۱۹	۵۳/۱۷

جدول (۱۴): نتایج ارزیابی الگوریتم SVM نمونه‌های آموزشی (رویکرد ترکیبی)

نماد	% دقت کلی	% فراخوان نظر مثبت	% فراخوان نظر منفی	ماتریس درهم ریختگی		
				خطا	کاهش	افزایش
اخابر	۷۲/۴۱	۶۸/۷۸	۷۵/۶۵	افزایش	۱۴۱	۳۱/۲۲
				کاهش	۵۶	۲۴/۳۵
حفاری	۷۲/۱۳	۷۰/۵۶	۷۳/۸۳	افزایش	۱۶۳	۲۹/۴۴
				کاهش	۵۶	۲۶/۱۷
خسپا	۷۳/۷۵	۷۲/۰۰	۷۵/۸۳	افزایش	۱۸۰	۲۸/۰۰
				کاهش	۵۱	۲۴/۱۷
خودرو	۷۰/۶۵	۶۷/۹۲	۷۲/۹۸	افزایش	۱۴۴	۳۲/۰۸
				کاهش	۶۷	۲۷/۰۲
ذوب	۷۳/۷۳	۷۲/۵۳	۷۵/۰۰	افزایش	۱۶۹	۲۷/۴۷
				کاهش	۵۵	۲۵/۰۰
ستران	۷۲/۶۴	۷۳/۴۴	۷۱/۶۵	افزایش	۱۷۷	۲۶/۵۶
				کاهش	۵۵	۲۸/۳۵

ماتریس درهم ریختگی			خطا	کاهشی	افزایش	٪فراخوان نظرمنفی	٪فراخوان نظرمثبت	٪دقت کلی	نماد
افزایشی	کاهشی	خطا							
افزایشی	کاهشی	خطا	۲۶/۸۸	۶۸	۱۸۵	۷۶/۶۳	۷۳/۱۲	۷۵/۲۱	شتران
کاهشی	افزایشی	خطا	۲۳/۳۷	۱۷۰	۴۹				
افزایشی	کاهشی	خطا	۲۸/۴۰	۶۹	۱۷۴	۷۶/۱۵	۷۱/۶۰	۷۳/۷۵	شبندر
کاهشی	افزایشی	خطا	۲۳/۸۵	۱۶۶	۵۲				
افزایشی	کاهشی	خطا	۳۲/۸۰	۶۱	۱۲۵	۷۲/۳۵	۶۷/۲۰	۶۹/۶۶	شپلی
کاهشی	افزایشی	خطا	۲۷/۶۵	۱۲۳	۴۷				
افزایشی	کاهشی	خطا	۲۵/۷۷	۶۷	۱۹۳	۷۹/۳۱	۷۴/۲۳	۷۶/۴۶	فاراک
کاهشی	افزایشی	خطا	۲۰/۶۹	۱۶۱	۴۲				
افزایشی	کاهشی	خطا	۲۷/۳۱	۵۹	۱۵۷	۷۸/۱۰	۷۲/۶۹	۷۵/۵۵	کیسون
کاهشی	افزایشی	خطا	۲۱/۹۰	۱۸۹	۵۳				
افزایشی	کاهشی	خطا	۲۹/۸۰	۵۹	۱۳۹	۷۹/۸۳	۷۰/۲۰	۶۸/۵۸	وبصادر
کاهشی	افزایشی	خطا	۲۰/۱۷	۱۹۰	۴۸				

دقت استفاده از الگوریتم ماشین بردار پشتیبان برای رویکرد ترکیبی پیش‌بینی جهت حرکت سهام در نمودار (۱) بیانگر افزایش دقت در فراخوان نظرات منفی است می‌توان نتیجه گرفت که هر وقت احساسات منفی باشد بیشترین حرکت سهام با احساسات منفی رخ می‌دهد در نتیجه قدرت پیش‌بینی حرکات منفی سهام بهتر می‌شود.

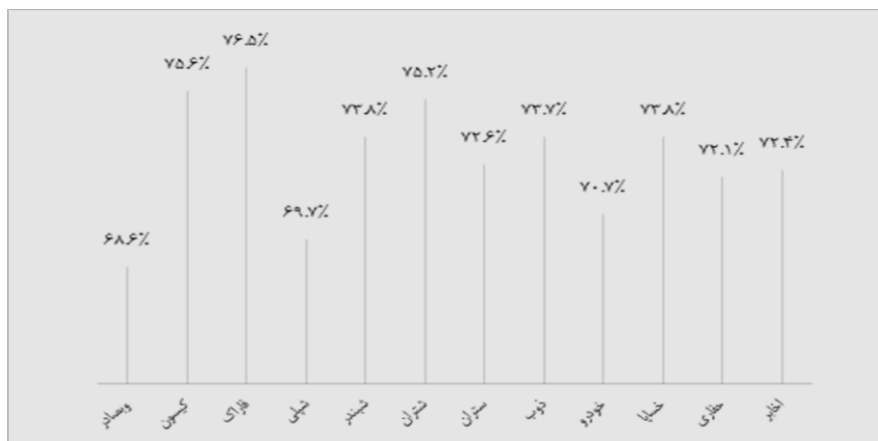


نمودار (۱) دقت فراخوان الگوریتم طبقه بندی SVM

۶- بحث و نتیجه گیری

در این پژوهش با استفاده از یک رویکرد ترکیبی تحلیل احساسات و نماگرهای تکنیکالی مدلی برای پیش‌بینی جهت حرکت سهام با سه نوع الگوریتم داده کاوی توسعه داده شد نتایج به دست آمده از الگوریتم ماشین بردار پشتیبان نشان داد دقت مدل در بخش تحلیل احساسات ۷۰/۲۰٪ و در بخش پیش‌بینی جهت حرکت سهام ۸۰/۳۹٪ بوده است. در پاسخ به سوالات پژوهش می‌توان گفت:

- با کمک احساسات و نماگرهای تکنیکالی می‌توان تا ۸۰٪ جهت حرکت سهام را پیش‌بینی نمود.
- طیف وسیعی از نظرات و احساسات معامله‌گران آنلاین و اثر آن بر جهت حرکت سهام بررسی شد در حالی که در پژوهش‌های گذشته فقط از اخبار و یا شایعات برای پیش‌بینی استفاده می‌شد
- بررسی رفتار معامله‌گران خرد در پیش‌بینی قیمت مهم است ولی اثرات نامتقارن سهامداران عمده می‌تواند در نتیجه مدل موثر باشد. در این پژوهش فقط رفتار معامله‌گران حقیقی و فعال در سهام یاب مورد بررسی قرار گرفت و فقدان بررسی رفتارهای حقوقی‌ها، سهامداران عمده و بازیگردان‌های سهام از نقاط ضعف این پژوهش به شمار می‌رود.
- با توجه به نمودار (۲) نتایج به دست آمده از الگوریتم SVM دقت پیش‌بینی جهت حرکت سهام را حدوداً ۸۰ درصد نشان می‌دهد که می‌توان برای توسعه مدل‌های آتی از آن استفاده نمود.



نمودار (۲) دقت الگوریتم SVM در پیش بینی جهت حرکت سهام

برای پژوهش‌های آتی می‌توان متغیرهای بنیادی موثر بر ارزش سهام، متغیرهای کلان اقتصادی نظیر نرخ بهره، بودجه و تغییر قیمت‌های جهانی و اثر آنها بر نظرات کاربران را نیز به عنوان متغیرهای ورودی به مدل اضافه کرد، استفاده از تحلیل گرافها و شبکه‌های پیوسته بین کاربران و و سطح دانش آنها در سرمایه‌گذاری و نیز توسعه یک واژه نامه مخصوص بازار سرمایه در زمینه تحلیل احساسات، مورد تحقیق و پژوهش‌های آتی قرار گیرد.

فهرست منابع

- * Fama, E. F. (1965). The behavior of stock-market prices. The journal of Business, Vol.38, no. 1, pp. 34-105.
- * Geweke, J. (1984). Inference and causality in economic time series models. Handbook of econometrics, Vol.2, pp.1101-1144.
- * Granger, C. W. (1969). Investigating causal relations by econometric models and cross-spectral methods. Econometrica: journal of the Econometric Society, pp.424-438.
- * J.Bollen, H. Mao, and X. Zeng. (2011). Twitter mood predicts the stock market. Journal of Computational Science.
- * Lee, H., Surdeanu, M., MacCartney, B., & Jurafsky, D. (2014, May). On the Importance of Text Analysis for Stock Price Prediction. In LREC (Vol. 2014, pp. 1170-1175).
- * Mittal, Anshul, and Arpit Goel. (2012). Stock Prediction Using Twitter Sentiment Analysis. Stanford University, CS229. Available online: <http://cs229.stanford.edu/proj2011/GoelMittalStockMarketPredictionUsingTwitterSenti>
- * Moshari, M., Didekhani, H., Dameghani, K.K., Abbasi, E. (2020). "Investigating the Predictability of Starting Point and Ending Short-Term Trend of Stock Price Using the Bayesian Likelihood Network", Journal of Financial Management Strategy, Vol. 8, No. 28 Spring

- * Schumaker, R. P., & Chen, H. (2009). Textual analysis of stock market prediction using breaking financial news: The AZF in text system. *ACM Transactions on Information Systems*, 27(2)
- * Vapnik, V. (1995). "The Nature of Statistical Learning Theory", New York, NY: Springer.
- * Vatanparast Mohammadreza, Asadi masoud, Mohammadi Shaban, Babaei abbas. 2019. Stock price prediction based on LM-BP neural network and over-point estimation by counting time intervals: Evidence from the Stock Exchange. *FINANCIAL summer*, Volume 10, Number 39; Page(s) 193 to 218.

یادداشت‌ها

- 1 Decision Tree
- 2 Support Vector Machine
- 3 Naive Bayes
- 4 Term Frequency - Inverse Document Frequency
- 5 Conditional probability
- 6 Vapnik et al
- 7 Granger Causality
- 8 Vector autoregressive
- 9 Geweke
- 10 Stationary
- 11 Dickey Fuller
- 12 Schumaker
- 13 Chen
- 14 Google-Profile of Mood States
- 15 Dow Jones Industrial Average
- 16 Mittal and Goel
- 17 Text mining