



دانشگاه صنعتی امیرکبیر  
(پلی تکنیک تهران)  
دانشکده صنایع

پایان نامه کارشناسی

معاملات الگوریتمی مالی ترکیبی با استفاده از الگوریتم‌های جنگل  
تصادفی، ماشین‌های بردار پشتیبان، شبکه‌های عصبی و نزدیک‌ترین  
همسایه

نگارش  
مسعود هادی نجف‌آبادی

استاد راهنما  
دکتر مسعود ماهوتچی

خرداد ماه ۱۴۰۰

صفحه فرم ارزیابی و تصویب پایان نامه - فرم تأیید اعضاء کمیته دفاع

اینجانب مسعود هادی نجفآبادی متعهد می‌شوم که مطالب مندرج در این پایان‌نامه حاصل کار پژوهشی اینجانب تحت نظارت و راهنمایی اساتید دانشگاه صنعتی امیرکبیر بوده و به دستاوردهای دیگران که در این پژوهش از آن‌ها استفاده شده است مطابق مقررات و روال متعارف ارجاع و در فهرست منابع و مآخذ ذکر گردیده است. این پایان‌نامه قبلاً برای احراز هیچ مدرک هم‌سطح یا بالاتر ارائه نگردیده است.

در صورت اثبات تخلف در هر زمان، مدرک تحصیلی صادر شده توسط دانشگاه از درجه اعتبار ساقط بوده و دانشگاه حق پیگیری قانونی خواهد داشت.

کلیه نتایج و حقوق حاصل از این پایان‌نامه متعلق به دانشگاه صنعتی امیرکبیر می‌باشد. هرگونه استفاده از نتایج علمی و عملی، واگذاری اطلاعات به دیگران یا چاپ و تکثیر، نسخه‌برداری، ترجمه و اقتباس از این پایان‌نامه بدون موافقت کتبی دانشگاه صنعتی امیرکبیر ممنوع است. نقل مطالب با ذکر مآخذ بلامانع است.

مسعود هادی نجفآبادی

امضا

## چکیده

روند تغییرات شاخص همواره به عنوان یکی از معیارهای سرمایه‌گذاری مدنظر قرار می‌گیرد. این پژوهش یک سیستم پشتیبان تصمیم برای معاملات الگوریتمی در بازارهای مالی ارائه می‌دهد که از رویکردی جدید در اتخاذ تصمیمات خودکار استفاده می‌کند. این رویکرد با ترکیب الگوریتم‌های جنگل تصادفی، ماشین بردار پشتیبان، شبکه عصبی و نزدیک‌ترین همسایه به پیش‌بینی سیگنال‌های خرید، نگهداری یا فروش می‌پردازد. داده‌های ورودی مدل، قیمت و حجم معاملات پانصد سهم شاخص S&P500 و قیمت سهام صندوق SPY به عنوان معیار شاخص S&P500 از ابتدای ماه ژانویه سال ۲۰۱۵ تا پایان ماه ژوئن سال ۲۰۲۱ بوده است. همچنین به جهت کاهش نویز ورودی به سیستم اندیکاتورهای میانگین متحرک وزن‌دار پنج روزه و شاخص قدرت نسبی نیز مورد استفاده قرار گرفته است. در مرحله نخست از الگوریتم جنگل تصادفی برای تعیین صد سهم دارای بیشترین اهمیت استفاده شده است، سپس با توجه به نتیجه‌ی مرحله‌ی اول هر یک از الگوریتم‌ها بر اساس داده‌های مربوط به ابتدای سال ۲۰۱۵ تا یک روز قبل از روز مورد پیش‌بینی، به تولید سیگنال خرید، نگهداری یا فروش می‌پردازند و سپس نتایج پیش‌بینی هر یک از الگوریتم‌ها با معیار دقت مدل سنجیده می‌شود. در مرحله‌ی پایانی با استفاده از تکنیک بگینگ<sup>۱</sup> که یک متالگوریتم برای ترکیب نتایج چند پیش‌بینی‌کننده بر اساس رای اکثریت است، استفاده می‌شود. به این نحو که نتایج پیش‌بینی الگوریتم‌های مرحله‌ی قبل در هر روز به عنوان متغیر ورودی در اختیار الگوریتم بگینگ قرار گرفته و این الگوریتم سیگنال نهایی را بر اساس رای اکثریت برای روز آتی تولید می‌کند.

## واژه‌های کلیدی:

معاملات الگوریتمی، بازار سهام، پیش‌بینی مالی، الگوریتم ترکیبی.

<sup>۱</sup> Bagging

صفحه	فهرست مطالب
۱.....	فصل اول مقدمه.....
۴.....	فصل دوم مبانی نظری و مروری بر ادبیات پژوهش.....
۵.....	۱-۲- شبکه عصبی (ANN).....
۷.....	۲-۲- ماشین بردار پشتیبان (SVM).....
۸.....	۳-۲- درخت تصمیم.....
۹.....	۴-۲- یادگیرنده‌های گند (یادگیری از طریق همسایه ها).....
۱۰.....	۱-۴-۲- دسته‌بندی‌های k همسایه‌ی نزدیک.....
۱۰.....	۵-۲- تکنیک‌هایی جهت بهبود صحت دسته‌بندی.....
۱۱.....	۱-۵-۲- روش bagging.....
۱۱.....	۲-۵-۲- روش boosting و الگوریتم AdaBoost.....
۱۱.....	۳-۵-۲- جنگل‌های تصادفی.....
۱۳.....	فصل سوم مروری بر پیشینه پژوهش.....
۱۸.....	فصل چهارم روش‌شناسی پژوهش.....
۱۹.....	۱-۴- تعریف مساله.....
۲۰.....	۲-۴- مراحل اجرای پژوهش.....
۲۰.....	۱-۲-۴- مرحله اول: شناسایی و جمع‌آوری متغیرهای مسئله.....
۲۰.....	۲-۲-۴- مرحله دوم: انتخاب صد متغیر دارای بیشترین اهمیت.....
۲۰.....	۳-۲-۴- مرحله سوم: پیش‌بینی و ارزیابی اولیه.....
۲۱.....	۴-۲-۴- مرحله چهارم: پیش‌بینی بر اساس داده‌های آزمایشی و ارزیابی مدل‌ها.....
۲۲.....	۵-۲-۴- مرحله پنجم: ترکیب نتایج پیش‌بینی الگوریتم‌ها و تولید سیگنال نهایی.....
۲۴.....	فصل پنجم بررسی نتایج.....
۲۵.....	۱-۵- جزئیات پیاده‌سازی مدل.....
۲۵.....	۲-۵- سنجش عملکرد مدل.....
۲۷.....	۳-۵- مقایسه بازده.....
۳۰.....	فصل ششم جمع‌بندی، نتیجه‌گیری و پیشنهادات.....
۳۱.....	۱-۶- جمع‌بندی و نتیجه‌گیری.....
۳۱.....	۲-۶- پیشنهادات.....
۳۳.....	منابع و مراجع.....

پیوست‌ها.....	۳۵
---------------	----

صفحه

فهرست اشکال

شکل ۴-۱: مدل پیش‌بینی.....	۲۳
شکل ۵-۱: مقایسه‌ی بازده مدل و بازده استراتژی خرید و نگهداری در بازه‌ی زمانی ابتدای ژانویه‌ی ۲۰۲۰ تا ابتدای جولای ۲۰۲۱.....	۲۷
شکل ۵-۲: مقایسه‌ی بازده مدل و بازده استراتژی خرید و نگهداری در بازه‌ی زمانی ابتدای ژانویه‌ی ۲۰۲۱ تا ابتدای جولای ۲۰۲۱.....	۲۸
شکل ۵-۳: سیگنال‌های ایجاد شده توسط مدل در بازه‌ی زمانی آزمون (ابتدای ژانویه‌ی ۲۰۲۰ تا ابتدای جولای ۲۰۲۱).....	۲۸
شکل ۵-۴: سیگنال‌های ایجاد شده توسط مدل در بازه‌ی زمانی آزمون (ابتدای ژانویه‌ی ۲۰۲۱ تا ابتدای جولای ۲۰۲۱).....	۲۹

صفحه

فهرست جداول

جدول ۵-۱: عملکرد الگوریتم‌ها در پیش‌بینی با استفاده از داده‌های آموزشی در بازه‌ی زمانی ۲۰۱۵ الی ۲۰۲۰	۲۶
جدول ۵-۲: عملکرد الگوریتم‌ها در پیش‌بینی با استفاده از داده‌های آموزشی در بازه‌ی زمانی ۲۰۱۵ الی ۲۰۲۱	۲۶
جدول ۵-۳: عملکرد الگوریتم‌ها در پیش‌بینی با استفاده از داده‌های آزمایشی در بازه‌ی زمانی ابتدای ۲۰۲۰ الی ابتدای جولای ۲۰۲۱	۲۶
جدول ۵-۴: عملکرد الگوریتم‌ها در پیش‌بینی با استفاده از داده‌های آزمایشی در بازه‌ی زمانی ابتدای ۲۰۲۱ الی ابتدای جولای ۲۰۲۱	۲۷
جدول پ-۱: شرح کد مدل ارائه شده	۳۵



## فهرست علائم و اختصارات

### علائم لاتین

جنگل تصادفی	RF
پیش‌بینی‌کننده‌ی چندلایه	MLP
ماشین بردار پشتیبان	SVM
K نزدیک‌ترین همسایه	KNN
میانگین متحرک وزن‌دار پنج روزه	WMA5
شاخص قدرت نسبی	RSI
دقت	ACC
مثبت صحیح	TP
منفی صحیح	TN
مثبت غلط	FP
منفی غلط	FN

## فصل اول

### مقدمه

یکی از مهم‌ترین دستاوردهای علم قابلیت پیش‌بینی بخشیدن به متغیرها و پدیده‌ها می‌باشد. پژوهشگران علوم مالی نیز با استفاده از ابزارهای مختلف به دنبال طراحی مدل‌هایی هستند که به وسیله آن متغیرها و حوادث مدنظرشان در بازارهای مالی را پیش‌بینی کنند. پیش‌بینی سود (جوگ و مک کنومی<sup>۱</sup>، ۲۰۰۳) پیش‌بینی ورشکستگی (وندا<sup>۲</sup>، ۲۰۰۴) و پیش‌بینی جریان نقدی (براجت<sup>۳</sup>، ۲۰۰۷) از این دست پژوهش‌هاست [۱].

به جرات می‌توان گفت که پیش‌بینی قیمت سهام از مهم‌ترین اهداف در علوم مالی و سرمایه‌گذاری است (هاموز<sup>۴</sup>، ۲۰۱۵). امروزه پیش‌بینی قیمت سهام نه تنها بسیار چالش برانگیز، بلکه بسیار مورد علاقه سرمایه‌گذاران می‌باشد. پیش‌بینی قیمت سهام از دو جنبه حائز اهمیت است. اول آنکه پیش‌بینی دقیق قیمت سهام برای بهینه‌سازی پرتفویهای سرمایه‌گذاری و اجرای راهبردهای سرمایه‌گذاری بسیار تأثیرگذار است (چانگ و چین<sup>۵</sup>، ۲۰۰۵). دوم آنکه با استفاده از روش‌های پیش‌بینی مطلوب می‌توان بازدهی را در سطح مشخصی از ریسک افزایش داد و بالعکس (کیومر<sup>۶</sup>، ۲۰۰۶) [۱].

در این میان شاخص به عنوان یک معیار آماری، قابلیت مقایسه وضعیت کنونی را نسبت به گذشته فراهم آورده و بررسی و تحلیل آن اطلاعات مفیدی را به کارشناسان و افراد ذی‌ربط در آن حوزه ارائه می‌دهد. شاخص‌های قیمت سهام در تمام بازارهای مالی دنیا، به مثابه یکی از مهم‌ترین معیارهای سنجش عملکرد بازار سهام، اهمیت زیادی دارند و شاید مهم‌ترین دلیل توجه روزافزون به آن‌ها، این نکته باشد که شاخص‌ها از تجمیع حرکت‌های قیمتی سهام تمام شرکت‌ها یا طبقه‌ی خاصی از شرکت‌های موجود در بازار به دست می‌آیند (ژیان ژو وانگ، ژنگ، گو و ژو ژی وانگ<sup>۷</sup>، ۲۰۱۱) [۲].

<sup>1</sup> Jog, V. & McConomy, B. J.

<sup>2</sup> Wallace Wanda, A.

<sup>3</sup> Brochet, F.

<sup>4</sup> Al-Hmouz, R.

<sup>5</sup> Cheung, Y. W., Chinn, M. D.

<sup>6</sup> Kumar, M.

<sup>7</sup> Wang, J. Z., Wang, J. J., Zhang, Z. G., & Guo, S. P.

داده‌های تاریخی نشان می‌دهد ویژگی‌های پیچیده شاخص کل قیمت، مانند غیرخطی بودن، عدم قطعیت، نوسان و پویایی، پیش‌بینی آن را دشوار می‌کند و نتایج پیش‌بینی را با عدم قطعیت زیادی مواجه می‌سازد که خود تاثیر شایان توجهی در بازده سرمایه‌گذاران، صندوق‌های سرمایه‌گذاری، نهادهای سرمایه‌گذاری و سایر فعالان این حوزه به همراه دارد (خسروی‌نژاد و شعبانی، ۱۳۹۳) [۲].

سه نوع تحلیل شناخته‌شده در بازارهای مالی مورد استفاده قرار می‌گیرد. تحلیل بنیادی بر پایه‌ی عملکرد شرکت‌ها و رشد سودآوری آن‌ها بنا شده است و تحلیل مالی-رفتاری حوزه‌ای از دانش مالی است که از نظریه‌های مبتنی بر روانشناسی برای توضیح رفتار بازارهای مالی بهره می‌گیرد. تحلیل تکنیکال سومین روش است که بر پایه سابقه معاملات یک دارایی مالی از طریق نمودار قیمت و فرمول‌های ریاضی که اندیکاتورهای تکنیکال نامیده می‌شوند، بنا شده است. در سال‌های اخیر از هوش مصنوعی نیز برای پیش‌بینی بازار استفاده شده است که ترکیب آن با تحلیل تکنیکال می‌تواند منجر به ایجاد سیستم‌های خودکار معاملاتی و الگوریتمی شود. از اولین کارها که با ترکیب تحلیل تکنیکال و الگوریتم‌های بهینه‌سازی و هوش مصنوعی سعی در ایجاد یک سیستم معاملاتی خودکار داشته است، می‌توان به اسکابار و کلوت<sup>۱</sup> (۲۰۰۲) اشاره کرد [۳].

در این پژوهش تلاش شده است تا با بهره‌گیری از اندیکاتورهای تحلیل تکنیکال در مرحله آغازین قابلیت پیش‌بینی‌پذیری را در اثر کاهش نویز ورودی به مدل افزایش داده و سپس با ارائه‌ی ترکیب جدیدی از الگوریتم‌های داده‌کاوی، در یک فرآیند سه مرحله‌ای به تولید سیگنال‌های معاملاتی برای روز آتی پرداخته شده است.

در مدل پیشنهاد شده در مرحله اول سهم‌های دارای بیشترین تاثیر انتخاب می‌شوند، سپس از چهار الگوریتم مختلف برای پیش‌بینی استفاده می‌شود و در نهایت با استفاده از یک الگوریتم چند لایه سیگنال‌های تولید شده در مرحله قبل ترکیب شده و سیگنال نهایی ارائه می‌شود.

بررسی‌های صورت گرفته نشان‌دهنده‌ی صحت قابل قبول سیگنال‌های تولید شده و بازده مناسب‌تر مدل از استراتژی خرید و نگهداری است.

<sup>1</sup> Skabar and Cloete

## فصل دوم

### مبانی نظری و مروری بر ادبیات پژوهش

در این قسمت، ما الگوریتم‌های دسته‌بندی با ناظر را که برای پیش‌بینی جهت حرکت قیمت استفاده شده است، مورد بحث و بررسی قرار می‌دهیم. دسته‌بندی، شکلی از تحلیل داده‌ها تلقی می‌شود که در آن مدل‌هایی جهت توصیف کلاس‌های مهمی از داده‌ها استخراج می‌شود. دسته‌بندی داده‌ها فرآیندی است که کار خود را در دو گام انجام می‌دهد: گام یادگیری (که در آن مدل ساخته می‌شود) و گام دسته‌بندی (که در آن جهت پیشگویی برچسب‌های کلاس از مدل ساخته شده در گام اول استفاده می‌شود). به دلیل آن که برچسب کلاس هر یک از تاپل‌های آموزشی مشخص شده‌اند، این گام همچنین به عنوان یادگیری با ناظر شناخته می‌شوند. در مقابل یادگیری بی‌ناظر (یا خوشه‌بندی) است، که در آن برچسب کلاس تاپل‌های آموزشی شناخته شده نیست و ممکن است تعداد یا مجموعه دسته‌هایی که در نهایت به دست خواهند آمد نیز از قبل مشخص نباشند [۴].

## ۲-۱- شبکه عصبی (ANN)

الگوریتم ANN از سیستم عصبی بیولوژیکی الهام گرفته شده است و می‌توان آن را به عنوان یک ارگانیسم کامل متشکل از تعداد زیادی واحد محاسباتی در نظر گرفت که برای حل مسئله با یکدیگر تعامل دارند. هر نورون سیگنال‌های سلول‌های عصبی همسایه را جمع‌آوری کرده و آن‌ها را به لایه بعدی منتقل می‌کند و در نتیجه سیگنال‌های "تحریک کننده" یا "بازدارنده" ایجاد می‌کند (گورونسکو<sup>۱</sup>، ۲۰۱۱). از این رو، هر نورون می‌تواند به عنوان یک پردازنده دیده شود که محاسبه ساده‌ای انجام می‌دهد؛ مانند تصمیم‌گیری در مورد ارسال کردن یا نکردن سیگنال به سلول‌های عصبی دیگر. یادگیری زمانی اتفاق می‌افتد که اثرات سیناپس‌ها تغییر کند، به عنوان مثال تأثیر یک نورون بر روی یک نورون دیگر تغییر می‌کند (تونچان<sup>۲</sup>، ۲۰۰۸).

مدل‌های ANN دارای سه ویژگی اصلی هستند:

- (۱) یک یا چند لایه از سلول‌های عصبی پنهان که بخشی از لایه‌های ورودی یا خروجی شبکه نیستند و امکان یادگیری را فراهم می‌کنند.

<sup>1</sup> Gorunescu

<sup>2</sup> Tunçhan

(۲) غیر خطی بودن قابل تمایز در فعالیت عصبی.

(۳) مدل اتصال شبکه از اتصال درجه بالایی برخوردار است.

این ویژگی‌ها همراه با یادگیری از طریق آموزش به حل مشکلات دشوار و متنوع کمک می‌کند. یادگیری از طریق آموزش در یک مدل ANN تحت نظارت را الگوریتم پس‌انتشار خطا<sup>۱</sup> نیز می‌نامند. الگوریتم BPN شبکه را براساس نمونه‌های ورودی-خروجی آموزش می‌دهد و یک سیگنال خطا پیدا می‌کند که تفاوت خروجی محاسبه شده و خروجی مورد نظر است. الگوریتم وزن سیناپسی گره‌های شبکه را متناسباً تنظیم می‌کند. بر اساس این اصل، یادگیری پس‌انتشار خطا در دو جهت رخ می‌دهد:

حرکت رو به جلو<sup>۲</sup>: در اینجا، یک ماتریس ورودی به شبکه ارائه می‌شود. هر ورودی به یک گره متصل می‌شود که یک سیگنال خروجی موقتی ایجاد می‌کند که در لایه بعدی توسط یک عملکرد انتقال به گره دیگری منتقل می‌شود. سیگنال ورودی از طریق لایه‌های شبکه به جلو انتشار می‌یابد و در انتهای خروجی شبکه به عنوان سیگنال خروجی ظاهر می‌شود. خروجی که در لایه خروجی محاسبه می‌شود با پاسخ مورد نظر مقایسه می‌شود و مقدار تفاوت، خطای آن گره را تعریف می‌کند. وزن‌های سیناپسی شبکه در طول حرکت رو به جلو ثابت می‌مانند.

حرکت رو به عقب<sup>۳</sup>: سیگنال خطایی که از نورون خروجی از آخرین لایه نشأت می‌گیرد، از طریق شبکه به عقب پخش می‌شود. این، شیب محلی هر نورون را در هر لایه محاسبه می‌کند و اجازه می‌دهد تا وزن سیناپسی شبکه مطابق با قانون دلتا تغییر کند. محاسبه بازگشتی با حرکت رو به جلو که توسط حرکت رو به عقب دنبال می‌شود، برای هر الگوی ورودی تا زمان همگرایی شبکه ادامه می‌یابد (گوا<sup>۴</sup>، ۲۰۰۷، مجومدر و هوسیان<sup>۵</sup>، ۲۰۰۷، مانتری و همکاران<sup>۶</sup>، ۲۰۱۴).

<sup>۱</sup> error back-propagation (BPN)

<sup>۲</sup> Forward Pass

<sup>۳</sup> Backward Pass

<sup>۴</sup> Quah

<sup>۵</sup> Majumder & Hussian

<sup>۶</sup> Mantri et al

BPN راه حل‌هایی را برای چندین مسئله خطی و غیر خطی مانند طبقه‌بندی، کنترل گیاه، پیش‌گویی، پیش‌بینی و علم رباتیک شناسایی می‌کند (مهرآرا و همکاران، ۲۰۱۰) [۵].

## ۲-۲- ماشین بردار پشتیبان (SVM)

SVM برای اولین بار توسط کورتس و وپنیک<sup>۱</sup> (۱۹۹۵) در زمینه تئوری یادگیری آماری و به حداقل رساندن ریسک ساختاری ارائه شده است. SVM در شناسایی الگو و مشکلات تخمین رگرسیون استفاده می‌شود و در حل مسائل تخمین وابستگی، پیش‌بینی و ساخت ماشین‌های هوشمند کاربرد دارد (اسمولا و شولکوف<sup>۲</sup>، ۲۰۰۴).

الگوریتم SVM به کمک یک نگاشت غیرخطی داده‌های آموزشی اولیه را به یک بعد بالاتر تبدیل می‌کند. در این بعد جدید به دنبال ابرصفحه‌ای<sup>۳</sup> بهینه می‌گردد، که تاپل‌های یک کلاس را از دیگر کلاس به صورت خطی تفکیک می‌کند. با یک نگاشت غیرخطی مناسب به یک بعد بالای کافی، داده‌های دو کلاس را همیشه می‌توان به کمک یک ابرصفحه تفکیک نمود. الگوریتم SVM این ابرصفحه را با کمک بردارهای پشتیبان<sup>۴</sup> (که اساساً تاپل‌های آموزشی هستند) و حاشیه‌ها<sup>۵</sup> (که با کمک بردارهای پشتیبان تعریف می‌شوند) پیدا می‌کند.

اگر داده‌ها به صورت خطی غیر قابل تفکیک باشند، یک طبقه‌بندی‌کننده SVM غیرخطی اعمال می‌شود. SVM بردارهای ورودی را با استفاده از یک تحول غیرخطی  $\Phi$  به یک فضای مشخصه چند بعدی تبدیل می‌کند و با استفاده از یک تابع هسته  $K(x, y)$  که در آن  $\{x_1 \dots x_n\}$  بردارهای ورودی و  $\{y_1 \dots y_n\}$  برچسب‌های آن‌ها  $\{1, -1\}$  هستند، در فضای مشخصه یک جداکننده خطی به عمل می‌آورد. هسته تابعی است که محصول نقطه‌ای  $(\phi(x), \phi(y))$  دو تصویر برداری در فضای مشخصه را

<sup>1</sup> Cortes & Vapnik

<sup>2</sup> Smola & Schölkopf

<sup>3</sup> Hyperplane

<sup>4</sup> Support Vectors

<sup>5</sup> Margine



برمی‌گرداند (جونگ و رجیا<sup>۱</sup>، ۲۰۰۸). اشکال مختلفی برای  $K(x, y)$  در ادبیات وجود دارد (نصرآبادی، ۲۰۰۷). در زیر چند نمونه آورده شده است.

خطی:

$$K(x, y) = x \cdot y$$

چند جمله‌ای:

$$K(x, y) = (x \cdot y)^d \text{ یا } K(x, y) = (1 + x \cdot y)^d$$

گوسی

$$K(x, y) = \exp \left[ -\frac{\|x - y\|^2}{2\delta^2} \right]$$

به دلیل آن که الگوریتم‌های SVM دارای قابلیت مدل‌سازی کران‌های تصمیم‌گیری غیرخطی پیچیده هستند، حتی سریع‌ترین آن‌ها نیز می‌تواند دارای سرعت پایینی در زمان آموزش باشد، اما صحت آن‌ها بسیار بالا است. در ضمن آن‌ها نسبت به دیگر روش‌ها کمتر دچار مشکل بیش‌برازش داده‌ها می‌شوند. بردارهای پشتیبان توصیف فشرده‌ای از مدل یادگیری‌شده را ارائه می‌دهند. از الگوریتم‌های SVM هم می‌توان برای پیشگویی عددی و هم برای دسته‌بندی داده‌ها استفاده کرد.

SVM همچنین با استفاده از رویکردهای "یکی در برابر یکی" یا "یکی علیه همه" برای مسائل چند طبقه‌ای<sup>۲</sup> اعمال می‌شود (مهرآرا و همکاران، ۲۰۱۰) [۵].

## ۲-۳- درخت تصمیم

یک درخت تصمیم همانطور که از نام آن مشخص است، یک ساختار درختی شبیه به فلوچارت دارد. هر گره داخلی (گره غیربرگ) در این درخت آزمونی را بر روی یک صفت خاصه نشان می‌دهد و هر شاخه نتیجه‌ی آزمون را نمایش می‌دهد و در هر گره برگ (یا گره پایانی) یک برچسب کلاس نگهداری می‌شود. برای ساخت درختان تصمیم به هیچ دانش خاص یا تنظیم پارامتری نیاز نیست. بنابراین برای یافتن اکتشافی دانش مناسب است. درختان تصمیم می‌توانند داده‌های چندبعدی را کنترل کنند. از نقطه نظر

<sup>۱</sup> Jung & Reggia

<sup>۲</sup> multi-class

بصری، هضم دانش ارائه شده در درختان تصمیم برای انسان راحت است. گام‌های دوگانه‌ی یادگیری و دسته‌بندی در استقراء درختان تصمیم ساده و سریع است و به طور کلی دارای صحت مناسبی هستند.

اواخر سال‌های ۱۹۷۰ و اوایل سال‌های ۱۹۸۰، محققان در حوزه‌ی یادگیری ماشین، الگوریتم ID3 را برای استقراء درخت تصمیم طراحی نمود. این کار تعمیم کارهای قبلی بر روی سیستم‌های یادگیری مفهوم<sup>۱</sup> تلقی می‌شد که توسط گروهی از محققین شرح داده شده بود. پس از این الگوریتم C4.5 (که نسخه‌ی بعدی ID3 محسوب می‌شود) ارائه شد. اغلب از این الگوریتم برای سنجش الگوریتم‌های جدیدتر در یادگیری با ناظر استفاده می‌شد. در سال ۱۹۸۴، گروهی از پژوهشگران در حوزه‌ی آمار کتابی را با نام درختان رگرسیون و دسته‌بندی<sup>۲</sup> (CART) منتشر کردند که در آن تولید درختان تصمیم دودویی بحث می‌شود. الگوریتم‌های ID3 و CART در زمان‌های تقریباً یکسانی و به صورت مستقل طراحی و پیشنهاد شدند و از رویکرد مشابهی برای یادگیری درختان تصمیم از تاپل‌های آموزش استفاده می‌کنند. این دو الگوریتم پایه و اساسی، باعث انجام سریع پژوهش‌هایی بر روی موضوع استقراء درخت تصمیم شدند [۴].

## ۲-۴- یادگیرنده‌های گند (یادگیری از طریق همسایه‌ها)

در رویکرد گند، یادگیرنده تا آخرین دقیقه انتظار می‌کشد و قبل از آن مدلی برای دسته‌بندی یک تاپل آزمایشی ساخته نمی‌شود. بنابراین در مواجهه با یک تاپل آموزشی یک یادگیرنده‌ی گند آن را ذخیره می‌کند (یا تنها پردازش کمی را انجام می‌دهد) و تا دریافت یک تاپل آزمایشی صبر می‌کند. تنها این یادگیرنده هنگامی که با تاپل آزمایشی روبه‌رو می‌شود، به تلاش برای دسته‌بندی تاپل مذکور بر اساس شباهت آن با تاپل‌های آموزشی ذخیره شده می‌پردازد. بر خلاف روش‌های یادگیری مشتاق، یادگیرنده‌های گند کار کمتری را هنگام آموزش و کار بیشتری در هنگام دسته‌بندی یا پیشگویی عددی انجام می‌دهند. ساختار داده‌ها در عملکرد آن‌ها تقریباً بی‌تاثیر است و به طور طبیعی از یادگیری افزایشی پشتیبانی می‌کند. آن‌ها قادر هستند فضای تصمیم پیچیده‌ای که شامل شکل‌هایی نظیر ابرچندضلعی‌ها می‌شود را مدل‌سازی کنند، که ممکن است این کار به سادگی توسط الگوریتم‌های دیگر یادگیری انجام نشود [۴].

<sup>۱</sup> Concept Learning Systems

<sup>۲</sup> Classification and Regression Trees

## ۲-۴-۱- دسته‌بندی‌های $k$ همسایه‌ی نزدیک<sup>۱</sup>

روش  $k$  نزدیک‌ترین همسایه برای اولین بار در اوایل سال‌های ۱۹۵۰ معرفی شد. این روش در مواجهه با داده‌های آموزشی حجیم با حجم عملیات بالایی روبه‌رو بود و به همین دلیل تا سال‌های ۱۹۶۰ که قدرت محاسبات کامپیوترها هنوز به اندازه‌ی کافی افزایش پیدا نکرده بود، از محبوبیت خوبی برخوردار نبود. از این الگوریتم به صورت گسترده‌ای در حوزه‌ی تشخیص الگو استفاده می‌شود.

دسته‌بندی‌های نزدیک‌ترین همسایه، یادگیری خود را بر اساس تشابه انجام می‌دهند، این کار با مقایسه‌ی تاپل آزمایشی و تاپل‌های آموزشی مشابه با آن صورت می‌گیرد. تاپل‌های آموزشی با کمک  $n$  صفت خاصه توصیف می‌شوند و هر تاپل در واقع نمایش نقطه‌ای در فضای  $n$  بعدی است. بدین ترتیب تمام تاپل‌های آموزشی در یک فضای  $n$  بعدی ذخیره می‌شوند. هرگاه با یک تاپل ناشناخته روبه‌رو می‌شوید، دسته‌بند  $k$  نزدیک‌ترین همسایه به دنبال  $k$  تاپل آموزشی است که شبیه‌ترین تاپل‌ها به تاپل ناشناخته هستند. این  $k$  تاپل آموزشی  $k$  همسایه‌ی نزدیک تاپل ناشناخته هستند.

دسته‌بندی‌های نزدیک‌ترین همسایه از مقایسه‌های مبتنی بر فاصله استفاده می‌کنند که در آن به طور طبیعی وزن هر یک از صفات خاصه برابر در نظر گرفته می‌شود. بنابراین هنگامی که داده‌ها حاوی نویز و صفات خاصه‌ی نامرتبط باشند، ممکن است با کاهش صحت دسته‌بند روبه‌رو شوید. به هر حال روش را می‌توان با وزن دهی صفات خاصه و حذف داده‌های نویز اصلاح نمود. انتخاب سنج‌های برای محاسبه‌ی فاصله می‌تواند امر مهمی باشد و ممکن است از سنج‌های دیگری به غیر از فاصله‌ی اقلیدسی برای این کار استفاده شود [۴].

## ۲-۵- تکنیک‌هایی جهت بهبود صحت دسته‌بندی

در این روش‌ها مدلی برای دسته‌بندی انتخاب می‌شود که ترکیبی از چندین دسته‌بند است. هر دسته‌بند رای خود را صادر می‌کند و نتیجه‌ی نهایی در مورد برچسب کلاس بر اساس این رای‌ها صادر می‌شود.

<sup>۱</sup> k-nearest-neighbor

## ۲-۵-۱- روش Bagging

مجموعه داده‌های  $D$  با تعداد  $d$  تاپل را در نظر بگیرید. روش bagging این گونه عمل می‌کند: در تکرار  $i$  ام الگوریتم  $(i=1,2,\dots,k)$  مجموعه آموزشی  $D_i$  با کمک روش نمونه‌گیری با جایگزینی از مجموعه داده‌های  $D$  انتخاب می‌شود. به دلیل آن که از نمونه‌گیری با جایگزینی استفاده می‌شود، ممکن است برخی از تاپل‌های موجود در  $D$  در مجموعه‌ی  $D_i$  قرار نگیرند، در حالی که برخی دیگر، بیش از یک بار انتخاب شوند. با کمک هر یک از مجموعه‌های آموزشی  $D_i$  مدل  $M_i$  تولید می‌شود. برای دسته‌بندی تاپل جدیدی مانند  $X$  هر یک از دسته‌بندی‌های  $M_i$  برچسب کلاس پیشنهادی خود را به عنوان یک رای ارائه می‌دهند. دسته‌بند تلفیقی  $M^*$  با شمارش آراء، رای اکثریت را برای برچسب کلاس تاپل  $X$  انتخاب می‌کند. از روش bagging می‌توان جهت پیش‌گویی مقادیر پیوسته نیز استفاده کرد. برای این کار کافی است میانگین مقادیر برگردانده شده توسط دسته‌بندها را محاسبه کنیم [۴].

## ۲-۵-۲- روش Boosting و الگوریتم AdaBoost

در روش boosting به هر یک از تاپل‌های آموزشی نیز وزنی تخصیص داده می‌شود. پس از ساخت دسته‌بند  $M_i$  وزن‌ها تغییر خواهند کرد تا دسته‌بند  $M_{i+1}$  که پس از  $M_i$  تولید می‌شود، توجه بیشتری را بر روی تاپل‌هایی که به درستی توسط  $M_i$  دسته‌بندی نشده‌اند، داشته باشد. دسته‌بند نهایی  $M^*$  رای نهایی را با ترکیب رای‌های هر یک از دسته‌بندی‌های پایه محاسبه می‌کند، جایی که وزن رای هر یک از این دسته‌بندها تابعی از صحت آن است. در واقع وزن یک تاپل، دشواری دسته‌بندی آن تاپل را منعکس می‌کند. به صورتی که وزن بالاتر نشان می‌دهد اغلب این تاپل به درستی دسته‌بندی نشده است [۴].

## ۲-۵-۳- جنگل‌های تصادفی

در این بخش به سراغ روشی دیگر از روش‌های تلفیقی به نام جنگل‌های تصادفی می‌رویم. تصور کنید دسته‌بندی‌های استفاده شده در روش تلفیقی همگی از نوع درخت تصمیم هستند؛ بدین ترتیب این مجموعه تشکیل یک جنگل را خواهند داد. هر یک از درختان تصمیم با استفاده از یک انتخاب تصادفی صفات خاصی موجود در هر گره جهت تعیین انشعاب ساخته می‌شوند. به عبارت دیگر هر درخت بر اساس مقادیر یک بردار تصادفی ساخته می‌شود. این مقادیر دارای توزیع یکسانی برای تمام درختان

موجود در جنگل هستند و به صورت مستقلی نمونه گیری می‌شوند. برای دسته‌بندی نیز هر درخت رای خود را صادر و نتیجه نهایی با رای اکثریت تعیین می‌شود [۴].

## فصل سوم

### مروری بر پیشینه پژوهش

پنج گروه از پژوهشگران معتقد هستند که پیش‌بینی قیمت سهام امکان‌پذیر نیست. اولین گروه کسانی هستند که به فرضیه بازار کارا اعتقاد دارند. در بازار کارای سرمایه، اعتقاد بر این است که قیمت سهام انعکاسی از تمام اطلاعات مربوط به آن سهم است و تغییرات قیمت سهام دارای الگوی خاص قابل پیش‌بینی نیست.

دومین گروه کسانی هستند که به علت موثر بودن عوامل متعدد بر تغییرات قیمت سهام و رفتار آشوب‌گونه و غیر خطی تغییرات قیمت سهام پیش‌بینی قیمت سهام را امری غیرممکن می‌دانند (منجمی و همکاران ۱۳۸۸). در حقیقت پراکندگی قیمت سهام تحت تاثیر عوامل کلان اقتصادی مانند نرخ بهره و نرخ تورم، وقایع سیاسی مانند جنگ و تهدیدات بین‌المللی و منطقه‌ای، وقایع اجتماعی مانند اعتصاب‌ها و آشوب‌ها و عوامل رفتاری و روانی سرمایه‌گذاران قرار دارد. اگر قیمت یا بازده در بازارهای مالی با دقت بالایی قابل پیش‌بینی باشد، سیستم‌های پیش‌بینی کننده‌ی قیمت تبدیل به یک دستگاه چاپ پول شده که ثروت زیادی را نصیب سرمایه‌گذاران می‌کند که این امر، در یک اقتصاد پایدار امکان‌پذیر نیست (گرنجر و تیمرمن<sup>۱</sup>، ۲۰۰۴). به بیان دیگر اگر تغییرات قیمت سهام در بازار قابل پیش‌بینی باشد (به طور تصادفی نباشد) احتمال وجود یک ثروت نامحدود برای سرمایه‌گذاران وجود خواهد داشت که با توجه به وضعیت حاضر و نتایج پژوهش‌های موجود این امکان وجود ندارد (گرانگر<sup>۲</sup>، ۱۹۹۱).

سومین گروه، نوع رویدادها در بازارهای مالی را دلیل اصلی عدم پیش‌بینی پذیری قیمت می‌دانند. هر رویدادی دارای دو ویژگی احتمال و شدت است. مدل‌های پیش‌بینی کننده قیمت، مبتنی بر احتمالات است و اساساً محققان به شدت رویدادها توجهی نمی‌کنند. طبق نظر این گروه به دلایل رفتاری، پیش‌بینی تغییرات قیمت با احتمال پایین ولی شدت بالا غیرقابل امکان است (نیکلاس تالب<sup>۳</sup>، ۲۰۱۲).

چهارمین گروه ماهیت بازارهای مالی را دلیل اصلی عدم پیش‌بینی پذیری قیمت سهام می‌دانستند. به نظر این گروه، ماهیت بازی به جمع صفر بودن در بازارهای مالی دلیل اصلی دشواری پیش‌بینی در این حوزه است. پیش‌بینی دقیق تر هوا باعث تغییر رفتار هوا نمی‌شود ولی پیش‌بینی دقیق قیمت سهام بر

<sup>1</sup> Timmermann, A., & Granger, C. W.

<sup>2</sup> Granger, C. W. J.

<sup>3</sup> Nassim Nicholas Taleb

خود قیمت سهام تاثیرگذار است. یک اقتصاددان با نام لوکاس چیزی را مطرح کرده بود که به نام خودش به عنوان نقد لوکاس معروف شد. نقد لوکاس اظهار داشت که شاید پیش گویی اقتصاددانان بر روند اقتصاد تاثیر بگذارد که اثر آن پیش گویی را خنثی سازد. فرض کنیم اقتصاددانان تورم را پیش بینی کنند، خزانه داری و بانک مرکزی در واکنش به گفته های ایشان وارد عمل شده و با اعمال سیاست های پولی و مالی تورم را پایین می آورند (تانگ و لین<sup>۱</sup>، ۲۰۰۷).

پنجمین گروه عامل مهمی به نام *اقتصاد/اطلاعات* را مطرح می سازند. به نظر این گروه شکل گیری اطلاعات نامتقارن عامل مهمی در پیش بینی ناپذیری در بازارهای مالی است. همیشه کسانی (مدیران ارشد شرکت ها) هستند که دارای اطلاعات محرمانه باشند (نیکواقبال و همکاران، ۱۳۹۲).

با این حال گروهی از پژوهشگران با توجه به نتایج حاصل از تحقیقات شان پیش بینی در بازارهای مالی و قیمت سهام را امری بسیار دشوار ولی امکان پذیر می دانند (مالکی<sup>۲</sup>، ۲۰۰۳). ماهیت پیچیده، تکاملی، داشتن خاصیت دینامیکی و غیرخطی که به دلیل تعامل حوادث و شرایط اقتصادی بوجود می آید و انتظارات غیرعقلانی سرمایه گذاران پیش بینی قیمت سهام را به امری دشوار مبدل می سازد نه امری محال (ناکاموری، ۲۰۰۵). کشف و بهبود الگوریتم های پیش بینی کننده و امکان انجام محاسبات پیچیده به وسیله رایانه ها، پژوهشگران را در این امر دشوار یاری می رساند [۱].

برای پیش بینی قیمت سهام راه های متفاوتی وجود دارد. یک راه کاهش پیچیدگی با استخراج بهترین صفات یا انتخاب از بین ویژگی ها می باشد [۶]. این روش با کاهش پیچیدگی به پیش بینی قیمت با صحت بیشتر کمک می کند. بررسی رابطه میان شرکت های هم صنعت نیز به ایجاد مدل های یادگیری ماشینی کمک می کند. در اینجا به بررسی کوتاهی پیرامون کاربرد شبکه های عصبی، ماشین بردار پشتیبان، درخت تصمیم و نزدیک ترین همسایه در پیش بینی قیمت سهام می پردازیم.

<sup>1</sup> Tang, C. F., & Lean, H. H.

<sup>2</sup> Malkiel, B. G.



برگس<sup>۱</sup> و همکاران (۲۰۰۰) از شبکه عصبی با دو لایه پنهان برای پیش‌بینی معاملات آتی یورو/دلار با بهره‌گیری از بالاترین قیمت، پایین‌ترین قیمت، قیمت باز گشایی و قیمت بسته شدن، استفاده کردند [۷].

تیلاکاراتن<sup>۲</sup> و همکاران (۲۰۰۷) با بهره‌گیری از الگوریتم شبکه‌های عصبی به پیش‌بینی سیگنال معاملاتی روز آتی شاخص سهام معمولی استرالیا<sup>۳</sup> با استفاده از داده روز فعلی در قیمت پایانی S&P 500، شاخص FTSE 100 (انگلیس) و شاخص CAC 40 (فرانسه) به عنوان ورودی پرداخته‌اند. نویسندگان دریافتند که شبکه‌های عصبی رو به جلو<sup>۴</sup> عملکرد بهتری نسبت به شبکه‌های احتمالی<sup>۵</sup> دارند [۸].

تاکور و کومار<sup>۶</sup> (۲۰۱۷) روشی جدید از ترکیب الگوریتم‌های ماشین بردار پشتیبان وزنی و جنگل تصادفی برای ایجاد سیگنال‌های خرید / نگهداری / فروش ارائه کرده‌اند به گونه‌ای که در ابتدا با استفاده از جنگل تصادفی زیرمجموعه‌ای بهینه از طیف گسترده‌ای از ابزارهای تکنیکال انتخاب می‌شود و سپس از الگوریتم ماشین بردار پشتیبان وزن دار جهت پیش‌بینی استفاده می‌شود [۹].

پژوهش فونته و همکاران<sup>۷</sup> (۲۰۰۶) با به کار بردن الگوریتم ژنتیک سعی در بهینه کردن پارامترهای اندیکاتورهای مختلف تکنیکال داشته‌است. در بازار داخلی در خصوص طراحی یک سیستم معاملات الگوریتمی با استفاده از اندیکاتورهای تکنیکال، می‌توان به کار دستپاک و رستگار (۱۳۹۴) اشاره کرد. در این پژوهش با استفاده از اطلاعات اندیکاتورهای تکنیکال به پیش‌بینی روند قیمت سهم و تعیین میزان خریدنی، فروختنی و یا نگهداشتنی بودن آن سهم پرداخته شده‌است. این مقاله که در بازار بورس تهران انجام شده‌است، حاکی از بهتر بودن سیستم معاملات الگوریتمی طراحی شده مبتنی بر تکنیکال نسبت به استراتژی خرید و نگهداری است. در پژوهشی دیگر، فلاح پور و حکیمیان (۱۳۹۵) به طراحی یک

<sup>1</sup> Burgess

<sup>2</sup> Tilakaratne

<sup>3</sup> Australian All Ordinary (AORD) index

<sup>4</sup> ANN feedforward

<sup>5</sup> probabilistic networks

<sup>6</sup> Manoj Thakur & Deepak Kumar

<sup>7</sup> de la Fuente, D., Garrido, A., Laviada, J. and Gómez, A.

سیستم معاملات الگوریتمی از نوع معاملات زوجی در بورس اوراق بهادار تهران پرداخته اند که مطابق ادعای نویسندگان بازدهی چشمگیری نسبت به بازدهی معمولی سهام در مدت مشابه دارد [۳].

فلاح پور و علی پور (۱۳۹۳) به پیش بینی شاخص کل سهام بورس اوراق بهادار تهران با استفاده از شبکه های عصبی موجکی پرداختند. ابتدا از تبدیل موجک گسسته برای نویزدایی داده ها در سری زمانی استفاده کردند؛ سپس به کمک شبکه های عصبی به پیش بینی شاخص سهام پرداختند. بر اساس نتایج، عملکرد شبکه عصبی موجکی سطح خطای کمتری نسبت به شبکه عصبی معمولی در پیش بینی شاخص سهام داشت [۲].

مشاری و همکاران (۱۳۹۶) با کمک الگوریتم ژنتیک به بهینه سازی متغیرهای پژوهش پرداخته و سپس مدلی جهت پیش بینی نقاط طلایی ارائه نموده اند [۱].

سزار و همکارش (۲۰۱۸) با استفاده از شبکه عصبی حلقوی، اطلاعات سری های زمانی مالی را به تصاویر دو بعدی تبدیل کرده و با استفاده از آن به تحلیل و تشخیص نقاط خرید، نگهداری و فروش پرداخته اند. در واقع با استفاده از ۱۵ نوع تحلیل تکنیکال (با پارامترها و فواصل زمانی متفاوت برای هر سهم) بر روی شاخص داوجونز یک تصویر دو بعدی ایجاد کرده و سپس با تحلیل تصاویر سیگنال خرید، نگهداری و فروش ایجاد کرده اند [۱۰].

## فصل چهارم

### روش‌شناسی پژوهش

## ۴-۱- تعریف مساله

این پژوهش از نظر رویکرد کاربردی بوده و مبتنی بر پژوهش‌های میدانی است. جامعه آماری آن، اطلاعات جمع آوری شده از قیمت ۵۰۰ سهام موجود در شاخص S&P500<sup>۱</sup> در بازه زمانی ابتدای سال ۲۰۱۵ تا پایان ماه ژوئن سال ۲۰۲۱ می‌باشد که به جهت ارزیابی بهتر مدل، در دو مرتبه به دو سری داده آموزش و آزمون مختلف دسته‌بندی شدند. این پژوهش درصدد است تا مدل پیش‌بینی جدیدی برای پیش‌بینی قیمت سهام صندوق SPY<sup>۲</sup> به عنوان معیار شاخص S&P500 ارائه نماید، به گونه‌ای که در ابتدا از میان پانصد سهم موجود در شاخص S&P500 که اطلاعات هر یک در بازه زمانی آموزش ورودی مدل به شمار می‌رود، صد سهام دارای بیشترین اهمیت در پیش‌بینی را به عنوان ورودی‌های اصلی انتخاب کرده و با استفاده از ترکیب چند الگوریتم، نتیجه مطلوب‌تری ارائه دهد. مدل ترکیبی مدنظر در ابتدا با استفاده از هر یک از الگوریتم‌های MLP، RF، SVM و KNN به پیش‌بینی تغییر قیمت در روز آتی پرداخته و در مرحله دوم نتایج به دست آمده از الگوریتم‌های مذکور را با استفاده از متالگوریتم بگینگ ترکیب نموده و بر اساس رای اکثریت الگوریتم‌ها، سیگنال نهایی برای روز آتی را تولید می‌کند. لازم به ذکر است که پارامترهای هریک از الگوریتم‌های مورد استفاده، با استفاده از الگوریتم اعتبارسنجی متقابل<sup>۳</sup> جست‌وجوی شبکه‌ای<sup>۴</sup> در کتابخانه Scikit learn زبان برنامه نویسی پایتون، تعیین شده است.

<sup>۱</sup> S&P 500 به عنوان یکی از معیارهای اصلی بازار سهام ایالات متحده عمل می‌کند و سلامت مالی و ثبات اقتصاد را نشان می‌دهد.

<sup>۲</sup> SPDR S&P 500 Trust ETF که با عنوان SPY ETF نیز شناخته می‌شود، یکی از محبوب‌ترین صندوق‌هایی است که هدف آن دنبال کردن شاخص S&P 500 است.

<sup>۳</sup> Cross Validation

<sup>۴</sup> Grid Search

## ۲-۴-۲- مراحل اجرای پژوهش

### ۲-۴-۱- مرحله اول: شناسایی و جمع‌آوری متغیرهای مسئله

اطلاعات قیمت پایانی<sup>۱</sup> و حجم معاملات روزانه مربوط به ۵۰۰ سهام موجود در شاخص S&P500 و قیمت نماد SPY در بازه زمانی ابتدای سال ۲۰۱۵ تا پایان ماه ژوئن سال ۲۰۲۱ از طریق کتابخانه مالی یاهو<sup>۲</sup> در زبان برنامه‌نویسی پایتون<sup>۳</sup> استخراج گردید و در راستای کاهش نویز ورودی به مدل، با استفاده از این اطلاعات اندیکاتورهای شاخص قدرت نسبی<sup>۴</sup> و میانگین متحرک وزنی پنج روزه<sup>۵</sup> نیز برای هر یک از نمادها در بازه‌ی زمانی ذکر شده محاسبه شد.

### ۲-۴-۲- مرحله دوم: انتخاب صد متغیر دارای بیشترین اهمیت

در این مرحله با استفاده از قابلیت تعیین اهمیت متغیرها<sup>۶</sup> در الگوریتم جنگل تصادفی، صد سهام دارای بیشترین اهمیت در پیش‌بینی سیگنال‌های معاملاتی انتخاب می‌شوند. سپس پایگاه داده ورودی مدل که شامل قیمت پایانی، حجم معاملات، شاخص قدرت نسبی و میانگین متحرک وزنی پنج روزه برای صد سهم انتخاب شده، است به عنوان مجموعه متغیرهای مورد استفاده در پیش‌بینی، تشکیل می‌شود.

### ۲-۴-۳- مرحله سوم: پیش‌بینی و ارزیابی اولیه

در این گام داده‌ها به دو دسته آموزش و آزمایش تقسیم بندی می‌شوند، با توجه به اجرای مدل در دو بازه آموزش و آزمون متفاوت، در مرتبه اول داده‌های مربوط به ابتدای سال ۲۰۱۵ تا ابتدای سال ۲۰۲۰

<sup>۱</sup> Adjusted Close

<sup>۲</sup> Yahoo Finance

<sup>۳</sup> Python

<sup>۴</sup> Relative Strength Index (RSI)

<sup>۵</sup> 5 Day Weighted Moving Average (WMA5)

<sup>۶</sup> Feature Importance

به عنوان داده آموزش<sup>۱</sup> و داده‌های مربوط به ابتدای ۲۰۲۰ به بعد به عنوان داده آزمایش<sup>۲</sup> در نظر گرفته می‌شود و در مرتبه دوم داده‌های مربوط به ابتدای سال ۲۰۱۵ تا ابتدای سال ۲۰۲۱ به عنوان داده آموزش<sup>۳</sup> و داده‌های مربوط به ابتدای ماه ژانویه سال ۲۰۲۰ تا پایان ماه ژوئن سال ۲۰۲۱ به عنوان داده آزمایش<sup>۴</sup> در نظر گرفته می‌شود. در ادامه هر یک از الگوریتم‌های SVM، RF، MLP و KNN توسط داده‌های آموزشی برای پیش‌بینی سیگنال روز آتی آموزش داده می‌شوند و نتایج عملکرد هریک از الگوریتم‌ها برای پیش‌بینی توسط داده‌های آموزشی، با استفاده از ماتریس اختلال<sup>۵</sup> و دو معیار دقت<sup>۶</sup> و اطمینان<sup>۷</sup> سنجیده خواهد شد.

$$Precision(positive) = \frac{TP}{TP + FP}$$

$$Precision(negative) = \frac{TN}{TN + FN}$$

$$Recall(positive) = \frac{TP}{TP + FN}$$

$$Recall(negative) = \frac{TN}{TN + FP}$$

#### ۴-۲-۴- مرحله چهارم: پیش‌بینی بر اساس داده‌های آزمایشی و ارزیابی مدل‌ها

در این مرحله هر یک از الگوریتم‌ها بر اساس داده‌های مربوط به ابتدای سال ۲۰۱۵ تا یک روز قبل از روز مورد پیش‌بینی، به تولید سیگنال خرید، نگهداری یا فروش می‌پردازند (همانند معامله‌گری که در

<sup>۱</sup> Train

<sup>۲</sup> Test

<sup>۳</sup> Train

<sup>۴</sup> Test

<sup>۵</sup> Confusion Matrix

<sup>۶</sup> Precision

<sup>۷</sup> Recall

پایان روز فعلی برای پیش‌بینی سیگنال روز بعد، اطلاعات مربوط به روز جاری را نیز در اختیار خواهد داشت). در پایان نتایج پیش‌بینی هر یک از الگوریتم‌ها با معیار دقت<sup>۱</sup> مدل سنجیده می‌شود.

$$Accuracy(ACC) = \frac{TP + TN}{TP + FP + TN + FN}$$

#### ۴-۲-۵- مرحله پنجم: ترکیب نتایج پیش‌بینی الگوریتم‌ها و تولید سیگنال نهایی

در این مرحله نتایج الگوریتم‌های مرحله قبل با استفاده از متالگوریتم<sup>۲</sup> بگینگ<sup>۳</sup>، که یک الگوریتم ترکیبی است، سیگنال نهایی برای روز آتی را بر اساس رای اکثریت تولید می‌کند.

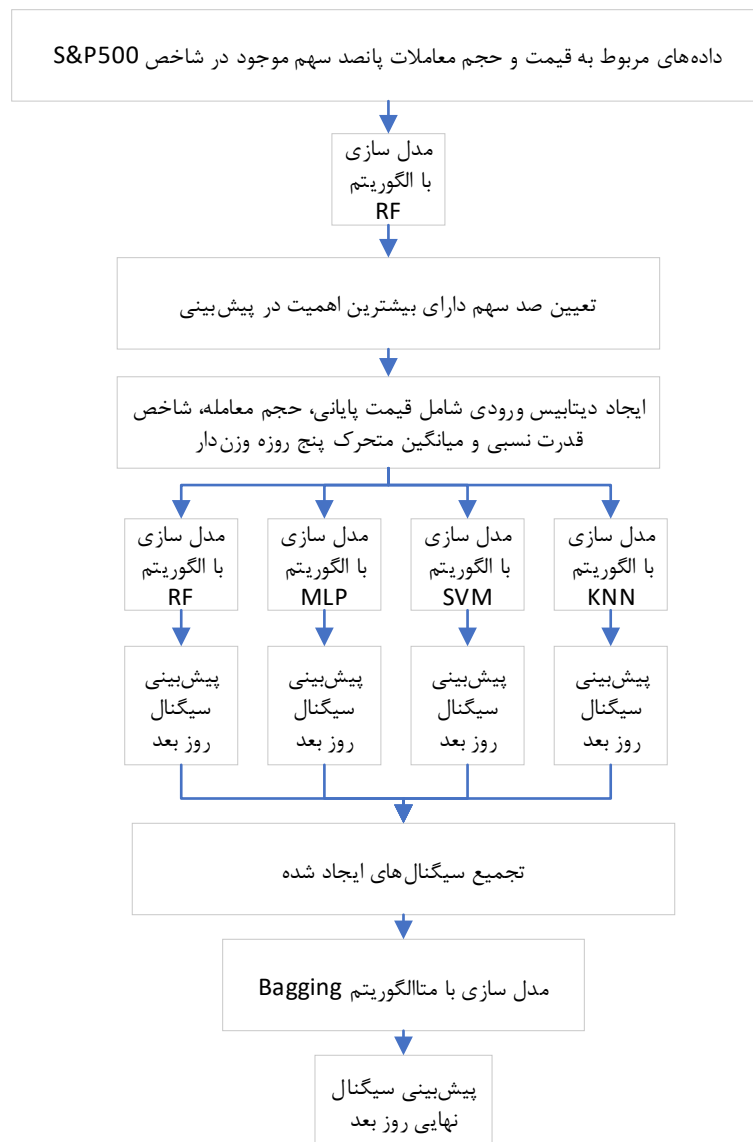
شکل ذیل مراحل اجرای مدل را نشان می‌دهد [شکل ۴-۱]:

---

<sup>۱</sup> Accuracy

<sup>۲</sup> Meta-Algorithm

<sup>۳</sup> Bagging (Bootstrap aggregating)



شکل ۴-۱: مدل پیش‌بینی



## فصل پنجم

### بررسی نتایج

## ۵-۱- جزئیات پیاده سازی مدل

تمامی الگوریتم‌ها از داده‌های مربوط به سهام شرکت‌های موجود در شاخص S&P500 به عنوان داده ورودی و از اطلاعات مربوط به قیمت سهام صندوق SPY به عنوان داده مورد پیش‌بینی و در بازه زمانی ابتدای سال ۲۰۱۵ الی ابتدای ماه جولای سال ۲۰۲۱ استفاده می‌کنند.

در این مطالعه با توجه به بررسی‌های موريس<sup>۱</sup> در رابطه با توزیع داده‌ها و مرزبندی تعیین سیگنال‌ها، فرض می‌شود که ۰.۵٪ افزایش (یا کاهش) در قیمت پلانی به اندازی کافی معقول است که حرکت مربوطه به عنوان سیگنال خرید (یا فروش) در نظر گرفته شود همچنین در صورتی که تغییر قیمت کمتر از ۰.۵٪ افزایش (یا کاهش) باشد، سیگنال نگهداری تولید خواهد شد [۵].

پارامترهای اصلی الگوریتم‌ها از جمله MLP و RF توسط قابلیت جستجوی شبکه‌ای<sup>۲</sup> کتابخانه سایکیت‌لرن<sup>۳</sup> و با اعتبارسنجی متقابل<sup>۴</sup> پنجگانه بیش از هزار و پانصد حالت مختلف تعیین خواهند شد.

## ۵-۲- سنجش عملکرد مدل

در مرحله آموزش هر یک از الگوریتم‌ها بر اساس دو معیار Precision و Recall مورد ارزیابی قرار گرفته‌اند. نتایج هر یک از الگوریتم‌ها در بازه‌های زمانی آموزش و آزمون متفاوت که در ادامه آورده شده است، حاکی از آن است که الگوریتم جنگل تصادفی دارای دقت صددرصد در پیش‌بینی بر اساس داده‌های آموزش است که این می‌تواند نشانه بیش‌برازش<sup>۵</sup> نیز باشد. همچنین واضح است که متالگوریتم بگینگ نیز خود را با نتایج الگوریتم جنگل تصادفی مطابقت داده و از دقت صددرصد در بازه زمانی آموزش برخوردار است.

<sup>۱</sup> S. A. Morris

<sup>۲</sup> Grid Search

<sup>۳</sup> Scikit Learn

<sup>۴</sup> Cross Validation

<sup>۵</sup> Over Fitting

جدول ۵-۱: عملکرد الگوریتم‌ها در پیش‌بینی با استفاده از داده‌های آموزشی در بازه‌ی زمانی

۲۰۱۵ الی ۲۰۲۰

Signal	MLP		RF		SVC		KNN		FINAL MODEL		Support
	Precision	Recall	Precision	Recall	Precision	Recall	Precision	Recall	Precision	Recall	
Sell	0.68	0.41	1.00	1.00	0.86	0.03	0.49	0.46	1.00	1.00	213
Hold	0.75	0.81	1.00	1.00	0.61	0.99	0.71	0.88	1.00	1.00	735
Buy	0.56	0.56	1.00	1.00	0.58	0.07	0.67	0.32	1.00	1.00	296
Accuracy	70%		100%		60%		67%		100%		

جدول ۵-۲: عملکرد الگوریتم‌ها در پیش‌بینی با استفاده از داده‌های آموزشی در بازه‌ی زمانی

۲۰۱۵ الی ۲۰۲۱

Signal	MLP		RF		SVC		KNN		FINAL MODEL		Support
	Precision	Recall	Precision	Recall	Precision	Recall	Precision	Recall	Precision	Recall	
Sell	0.97	0.82	1.00	1.00	0.67	0.01	0.49	0.52	1.00	1.00	284
Hold	0.87	1.00	1.00	1.00	0.59	0.97	0.70	0.84	1.00	1.00	818
Buy	0.95	0.79	1.00	1.00	0.50	0.17	0.58	0.32	1.00	1.00	396
Accuracy	90%		100%		57%		64%		100%		

در مرحله آزمون نمره دقت هر یک از الگوریتم‌ها، بر اساس درستی یا نادرستی پیش‌بینی سیگنال‌های روزانه، در دو مرحله و برای بازه‌های زمانی آموزش و آزمون متفاوت محاسبه شده است. نتایج حاکی از آن است که صحت پیش‌بینی توسط متالگوریتم بگینگ که در انتهای مدل قرار دارد، از صحت پیش‌بینی هریک از الگوریتم‌های دیگر بیشتر بوده و در نتیجه سیگنال‌های ایجاد شده توسط آن از صحت بالاتری برخوردار است.

جدول ۵-۳: عملکرد الگوریتم‌ها در پیش‌بینی با استفاده از داده‌های آزمایشی در بازه‌ی زمانی

ابتدای ۲۰۲۰ الی ابتدای جولای ۲۰۲۱

Signal	MLP		RF		SVC		KNN		FINAL MODEL		Support
	Precision	Recall	Precision	Recall	Precision	Recall	Precision	Recall	Precision	Recall	
Sell	0.22	0.25	0.30	0.24	0.35	0.06	0.26	0.35	0.31	0.27	93
Hold	0.43	0.38	0.51	0.66	0.43	0.92	0.52	0.60	0.52	0.66	146
Buy	0.30	0.31	0.41	0.34	0.46	0.17	0.41	0.25	0.42	0.34	138
Accuracy	32%		44%		43%		41%		45%		

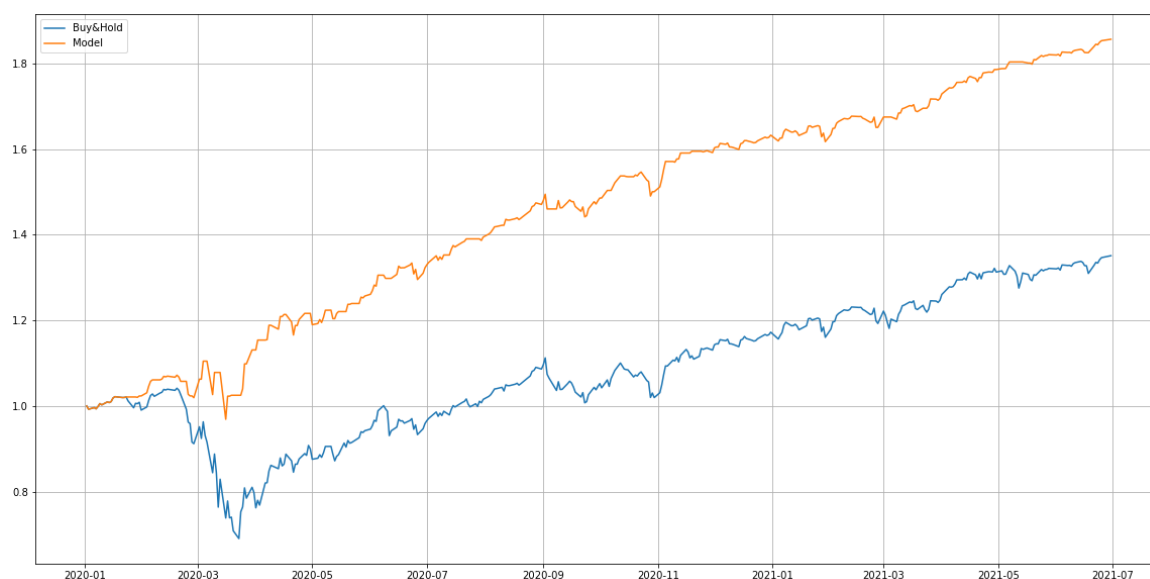
جدول ۴-۵: عملکرد الگوریتم‌ها در پیش‌بینی با استفاده از داده‌های آزمایشی در بازه‌ی زمانی

ابتدای ۲۰۲۱ الی ابتدای جولای ۲۰۲۱

Signal	MLP		RF		SVC		KNN		FINAL MODEL		Support
	Precision	Recall	Precision	Recall	Precision	Recall	Precision	Recall	Precision	Recall	
Sell	0.00	0.00	0.43	0.13	0.00	0.00	0.25	0.26	0.44	0.15	23
Hold	0.52	0.94	0.57	0.92	0.55	0.98	0.64	0.67	0.59	0.93	63
Buy	0.40	0.11	0.53	0.21	0.67	0.21	0.35	0.32	0.53	0.21	38
Accuracy	50%		56%		56%		48%		57%		

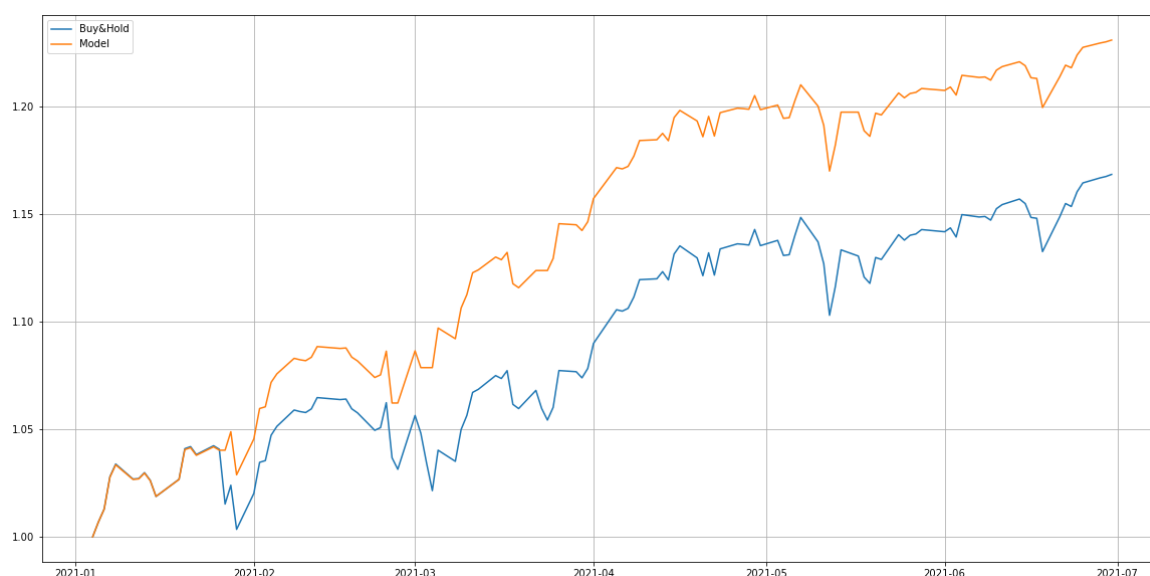
### ۵-۳- مقایسه بازده

نمودارهای ذیل بازده مدل را در مقایسه با بازده استراتژی خرید و نگهداری برای داده‌های آزمون و در دو بازه زمانی متفاوت نمایش می‌دهد که حاکی از عملکرد بهتر مدل می‌باشد.



شکل ۵-۱: مقایسه‌ی بازده مدل و بازده استراتژی خرید و نگهداری در بازه‌ی زمانی ابتدای ژانویه ۲۰۲۰ تا

ابتدای جولای ۲۰۲۱

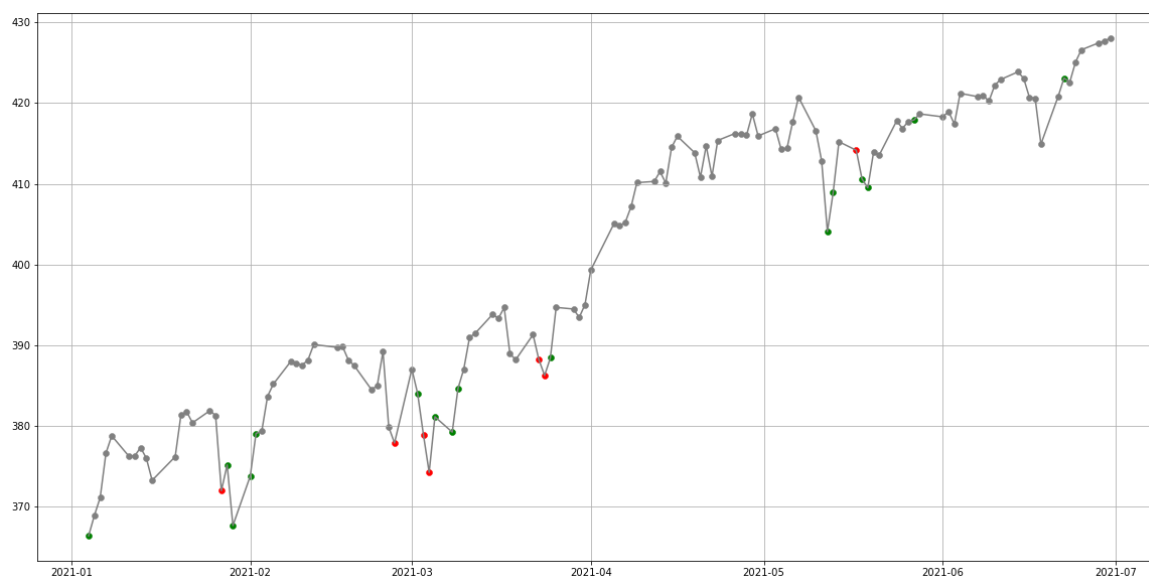


شکل ۵-۲: مقایسه‌ی بازده مدل و بازده استراتژی خرید و نگهداری در بازه‌ی زمانی ابتدای ژانویه ۲۰۲۱ تا ابتدای جولای ۲۰۲۱

نمودارهای ذیل سیگنال‌های تولید شده توسط مدل را بر روی نمودار قیمت سهام SPY و در دو بازه زمانی آزمون متفاوت نمایش می‌دهد. قابل ذکر است که سیگنال خرید با رنگ سبز، سیگنال فروش با رنگ قرمز و سیگنال نگهداری با رنگ خاکستری نمایش داده شده است.



شکل ۵-۳: سیگنال‌های ایجاد شده توسط مدل در بازه‌ی زمانی آزمون (ابتدای ژانویه ۲۰۲۰ تا ابتدای جولای ۲۰۲۱)



شکل ۴-۵: سیگنال‌های ایجاد شده توسط مدل در بازه‌ی زمانی آزمون (ابتدای ژانویه‌ی ۲۰۲۱ تا ابتدای جولای ۲۰۲۱)

بررسی سیگنال‌های تولید شده توسط مدل، بیانگر عملکرد مناسب مدل در مواجهه با تلاطمات قیمت می‌باشد.

## فصل ششم

### جمع‌بندی، نتیجه‌گیری و پیشنهادات

## ۶-۱- جمع‌بندی و نتیجه‌گیری

هدف اصلی این پژوهش دستیابی به نرخ سود بالاتر در بازار سرمایه است و همانطور که اشاره شد، تلاش شده است بر اساس تغییرات قیمت سهم‌های زیر مجموعه شاخص S&P500 تغییرات شاخص در روز معاملاتی بعد پیش‌بینی شود. مدل ارائه شده در این پژوهش در سه مرحله و با استفاده از الگوریتم‌های جنگل تصادفی، شبکه عصبی، ماشین بردار پشتیبان و نزدیکترین همسایه به پیش‌بینی سیگنال خرید، نگهداری و یا فروش می‌پردازد. به این صورت که در ابتدا با استفاده از الگوریتم جنگل تصادفی صد سهم دارای بیشترین اهمیت در پیش‌بینی انتخاب می‌شوند، سپس هر یک از الگوریتم‌های ذکر شده به پیش‌بینی سیگنال معاملاتی روز آتی می‌پردازند و در نهایت با استفاده از الگوریتم شبکه عصبی سیگنال نهایی بر پایه سیگنال‌های تولید شده در مرحله قبل ایجاد می‌شود.

صحت پیش‌بینی توسط الگوریتم نهایی از صحت پیش‌بینی هریک از الگوریتم‌های دیگر بیشتر بوده و میزان بازده ناشی از معامله بر اساس سیگنال‌های روزانه تولید شده توسط مدل برای داده‌های آزمون بیشتر از سود از ناشی استراتژی خرید و نگهداری است که بیانگر عملکرد مناسب مدل پیشنهادی است.

## ۶-۲- پیشنهادات

پیشنهاد می‌شود برای توسعه‌ی این روش از داده‌های تکمیلی همچون دیگر شاخص‌های بازار سهام ایالات متحده، شاخص بازار سهام دیگر کشورها و همچنین اندیکاتورهای دیگر استفاده شود. همچنین می‌توان با استفاده از الگوریتم‌های داده کاوی نوین و استفاده از روش‌های دیگر در جهت تعیین پارامترهای الگوریتم‌ها، عملکرد مدل را بهبود بخشید. همچنین به‌کارگیری متالگوریتم Boosting به عنوان الگوریتم نهایی می‌تواند مورد بررسی قرار گیرد.





## منابع و مراجع

- [۱] مشاری، محمد؛ دیده خانی، حسین؛ خلیلی دامغانی، کاوه؛ عباسی، ابراهیم؛ "طراحی مدل هوشمند ترکیبی جهت پیش‌بینی نقاط طلایی قیمت سهام"، فصلنامه علمی پژوهشی دانش سرمایه‌گذاری، سال هشتم، شماره ۲۹، صفحات ۴۵-۶۶، بهار ۱۳۹۸.
- [۲] درودی، دیاکو؛ ابراهیمی، سید بلبک؛ "ارائه‌ی روش هیبریدی نوین برای پیش‌بینی شاخص کل قیمت بورس اوراق بهادار"، تحقیقات مالی، دوره ۱۸، شماره ۴، صفحات ۶۳۲-۶۱۳، زمستان ۱۳۹۵.
- [۳] رستگار، محمدعلی؛ صداقتی پور، امین؛ "ارائه سیستم معاملات الگوریتمی برای قرارداد آتی سکه طلا مبتنی بر داده‌های درون-روزی"، فصلنامه علمی پژوهشی دانش سرمایه‌گذاری، سال هفتم، شماره ۲۸، صفحات ۴۹-۶۷، زمستان ۱۳۹۷.
- [۴] هان، ژیاوی؛ کمبر، میشلین؛ پی، ژان؛ اسماعیلی، مهدی؛ داده‌کاوی (مفاهیم و تکنیک‌ها)، انتشارات نیازدانش، تهران، ۱۳۹۳.
- [5] ERKARTAL, Bugra & Ozdamar, Linet. (2018). Generating Buy/Sell Signals for an Equity Share Using Machine Learning. Eurasian Journal of Business and Economics. 11. 85-105. 10.17015/ejbe.2018.022.04.
- [6] Mantri, J. K., Gahan, P., & Nayak, B. B. "Artificial neural networks--an application to stock market volatility. Soft-Computing in Capital Market: Research and Methods of Computational Finance for Measuring Risk of Financial Instruments", 179, 2014.
- [7] Burgess, A. N., & others. "A computational methodology for modelling the dynamics of statistical arbitrage", PHD Thesis: University of London, 2000.
- [8] Tilakaratne, C. D., Morris, S. A., Mammadov, M. A., & Hurst, C. P. "Predicting stock market index trading signals using neural networks", In Proceedings of the 14th Annual Global Finance Conference (GFC'07), Pages 171–179, 2007.
- [9] Manoj Thakur, Deepak Kumar. "A hybrid financial trading support system using multi-category classifiers and random forest", Applied Soft Computing, Volume 67, Pages 337-349, 2017.

- [10] Omer Berat Sezer, Ahmet Murat Ozbayoglu. “Algorithmic financial trading with deep convolutional neural networks: Time series to image conversion approach”, Applied Soft Computing, Volume 70, Pages 525-538, 2018.

## پیوست‌ها

جدول پ-۱: شرح کد مدل ارائه شده

```
#Feature Importance

y = spy['Signal'][spy.index.isin(data['Adj Close'].index)]
X = data['Adj Close'][data['Adj Close'].index.isin(spy.index)]

x_train = X.loc[X.index < '2020-01-01']
y_train = y.loc[y.index < '2020-01-01']

from sklearn.ensemble import RandomForestClassifier
fi_forest = RandomForestClassifier()
fi_forest.fit(x_train, y_train)
FI = pd.DataFrame()

FI['Stock'] = data['Adj Close'].columns
FI['Importance'] = fi_forest.feature_importances_

FI = FI.sort_values('Importance',ascending=False)
FI = FI.reset_index()
VIP = FI['Stock'].loc[:100].values
print (FI['Stock'].loc[:100].values)

# Change daily price to WMA5
df_WMA5 = pd.DataFrame()
for stock in data['Adj Close'].columns :
    df_WMA5[stock] = ta.WMA(data['Adj Close'][stock],timeperiod = 5)
df_WMA5.dropna(inplace=True)

# Change daily price to RSI
df_RSI = pd.DataFrame()
for stock in data['Adj Close'].columns :
    df_RSI[stock] = ta.RSI(data['Adj Close'][stock],timeperiod = 14)
df_RSI.dropna(inplace=True)

# Make a new dataset with WMA5 , Vol , RSI
df_data = pd.DataFrame()
df_data = pd.merge(df_WMA5[VIP]
                  ,data['Volume'][VIP]
                  ,right_index=True,left_index=True,suffixes=('_wma','_vol'))
df_data = pd.merge(df_data
                  ,df_RSI[VIP]
                  ,right_index=True,left_index=True,suffixes=('', '_rsi'))
df_data = pd.merge(df_data
                  ,spy[['Adj Close','Volume']]
                  ,right_index=True,left_index=True,suffixes=('', '_spy'))
```

```

df_data.dropna(axis=0,inplace=True)

#Make a Function to measure performance

from sklearn.model_selection import cross_val_score,
cross_val_predict
from sklearn.metrics import accuracy_score, classification_report,
confusion_matrix

def print_score(clf, X_train, y_train, X_test, y_test, train=True):
    '''
        print the accuracy score, classification report and confusion
        matrix of classifier
    '''
    if train:
        '''
            training performance
        '''
        print("Train Result:\n")
        print("accuracy score:
{0:.4f}\n".format(accuracy_score(y_train, clf.predict(X_train))))
        print("Classification Report: \n
{}\n".format(classification_report(y_train, clf.predict(X_train))))
        print("Confusion Matrix: \n
{}\n".format(confusion_matrix(y_train, clf.predict(X_train))))
        """
        res = cross_val_score(clf, X_train, y_train, cv=10,
scoring='accuracy')
        print("Average Accuracy: \t {0:.4f}".format(np.mean(res)))
        print("Accuracy SD: \t\t {0:.4f}".format(np.std(res)))
        print('\n')
        """
    elif train==False:
        '''
            test performance
        '''
        print("Test Result:\n")
        print("accuracy score:
{0:.4f}\n".format(accuracy_score(y_test, clf.predict(X_test))))
        print("Classification Report: \n
{}\n".format(classification_report(y_test, clf.predict(X_test))))
        print("Confusion Matrix: \n
{}\n".format(confusion_matrix(y_test, clf.predict(X_test))))

from sklearn.neural_network import MLPClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.svm import SVC
from sklearn.neighbors import KNeighborsClassifier
from sklearn.ensemble import AdaBoostClassifier
from sklearn.preprocessing import normalize

test_date = '2020-01-01'
y = spy['Signal'][spy.index.isin(df_data.index)]
X = df_data[df_data.index.isin(spy.index)]

result = pd.DataFrame()
result['Real'] = spy['Signal'][spy.index.isin(df_data.index)]
result['MLP'] = 0

```

```

result['RF'] = 0
result['SVM'] = 0
result['KNN'] = 0

#Train Part

#x_train = normalize(X.loc[X.index < test_date],axis=1)
x_train = X.loc[X.index < test_date]
y_train = y.loc[y.index < test_date]
#x_test = normalize(X.loc[X.index > test_date],axis=1)
x_test = X.loc[X.index > test_date]
y_test = y.loc[y.index > test_date]

mlp = MLPClassifier(hidden_layer_sizes = (10,),
                    activation = 'identity',
                    solver = 'lbfgs',
                    alpha = 0.001)
mlp.fit(x_train, y_train)
result['MLP'].loc[result.index < test_date] = mlp.predict(x_train)

forest = RandomForestClassifier(max_depth = 5,
                               max_features = 'log2',
                               min_samples_leaf = 10,
                               min_samples_split = 7,
                               n_estimators = 20)
forest = RandomForestClassifier()
ada_rf = AdaBoostClassifier(base_estimator=forest, n_estimators=100,
learning_rate=0.5)
ada_rf.fit(x_train, y_train.ravel())
#forest.fit(x_train, y_train)
result['RF'].loc[result.index < test_date] = ada_rf.predict(x_train)

svc = SVC()
svc.fit(x_train, y_train)
result['SVM'].loc[result.index< test_date] = svc.predict(x_train)

knn = KNeighborsClassifier()
knn.fit(x_train, y_train)
result['KNN'].loc[result.index< test_date] = knn.predict(x_train)

print('MLP :')
print_score(mlp, x_train, y_train, x_test, y_test, train=True)
print_score(mlp, x_train, y_train, x_test, y_test, train=False)
print('RF with Ada :')
print_score(ada_rf, x_train, y_train,x_test, y_test, train=True)
print_score(ada_rf, x_train, y_train,x_test, y_test, train=False)
print('SVC :')
print_score(svc, x_train, y_train,x_test, y_test, train=True)
print_score(svc, x_train, y_train,x_test, y_test, train=False)
print('KNN :')
print_score(knn, x_train, y_train,x_test, y_test, train=True)
print_score(knn, x_train, y_train,x_test, y_test, train=False)

# Test Part

```



```

grid_search.fit(x_train, y_train)

print(grid_search.best_score_)

grid_search.best_estimator_.get_params()

#GridSearch for MLP

from sklearn.model_selection import GridSearchCV

x_train =X.loc[X.index < test_date]
y_train = y.loc[y.index < test_date]

mlp_grid = MLPClassifier()

params_grid =
{'hidden_layer_sizes':[(10,),(100,),(20,),(200,),(50,),(500,)],
  'activation':['identity', 'logistic', 'tanh',
  'relu'],
  'solver':['lbfgs', 'sgd', 'adam'],
  'alpha':[0.001,0.0001,0.00001]}

grid_search = GridSearchCV(mlp_grid, params_grid,
                           n_jobs=-1, cv=5,
                           verbose=1, scoring='accuracy')

grid_search.fit(x_train, y_train)

print(grid_search.best_score_)

grid_search.best_estimator_.get_params()

def CLF_Report (Model,y_test,y_model):
    print(Model)
    print("accuracy score:
{0:.4f}\n".format(accuracy_score(y_test,y_model)))
    print("Classification Report: \n
{}\n".format(classification_report(y_test,y_model)))
    print("Confusion Matrix: \n
{}\n".format(confusion_matrix(y_test,y_model)))

y_test = result['Real'].loc[y.index > test_date]
rf_test = result['RF'].loc[y.index > test_date]
mlp_test = result['MLP'].loc[y.index > test_date]
svc_test = result['SVM'].loc[y.index > test_date]
knn_test = result['KNN'].loc[y.index > test_date]

CLF_Report('RF',y_test, rf_test)
CLF_Report('MLP',y_test, mlp_test)
CLF_Report('SVM',y_test, svc_test)
CLF_Report('KNN',y_test, knn_test)

result['past Real'] = result['Real'].shift(1)
result.dropna(axis=0,inplace=True)

#Final model with MLP
from sklearn.neighbors import KNeighborsClassifier

```



```

from sklearn.neural_network import MLPClassifier

#Train
y = result['Real']
X = result[['RF', 'MLP', 'SVM', 'past Real']]

x_train = X.loc[X.index < test_date]
y_train = y.loc[y.index < test_date]
x_test = X.loc[X.index > test_date]
y_test = y.loc[y.index > test_date]

result['Final'] = 0

#knn = KNeighborsClassifier()
#knn.fit(x_train, y_train)
#result['Final'].loc[result.index < test_date] =
knn.predict(x_train)

#mlp_clf = MLPClassifier()
#mlp_clf.fit(x_train, y_train)
#result['Final'].loc[result.index < test_date] =
mlp_clf.predict(x_train)

forest_clf = RandomForestClassifier()
ada_rf = AdaBoostClassifier(base_estimator=forest_clf,
n_estimators=10, learning_rate=0.5)
ada_rf.fit(x_train, y_train.ravel())
result['Final'].loc[result.index < test_date] =
ada_rf.predict(x_train)

y_train = result['Real'].loc[y.index < test_date]
final_train = result['Final'].loc[y.index < test_date]
CLF_Report('Final Train', y_train, final_train)

#Test
for i in y.loc[y.index >= test_date].index :
    x_train = X.loc[X.index < i].values
    y_train = y.loc[y.index < i]
    x_test = X.loc[X.index == i].values
    y_test = y.loc[y.index == i]

    #knn = KNeighborsClassifier()
    #knn.fit(x_train, y_train)
    #result['Final'].loc[result.index == i] =
    knn.predict(x_test)[0]

    #mlp_clf = MLPClassifier()
    #mlp_clf.fit(x_train, y_train)
    #result['Final'].loc[result.index == i] =
    mlp_clf.predict(x_test)[0]

    forest_clf = RandomForestClassifier()
    ada_rf = AdaBoostClassifier(base_estimator=forest_clf,
n_estimators=10, learning_rate=0.5)
    ada_rf.fit(x_train, y_train.ravel())
    result['Final'].loc[result.index == i] =
    ada_rf.predict(x_test)[0]

```

```

y_test = result['Real'].loc[y.index > test_date]
final_test = result['Final'].loc[y.index > test_date]
CLF_Report('Final Test : ' , y_test , final_test)

x_test = normalize(X.loc[X.index > test_date],axis=0)
y_test = y.loc[y.index > test_date]

Return = pd.DataFrame()
Return['Return'] = spy['Return'].loc[y_test.index]
Return['Buy&Hold'] = (1 + Return['Return']).cumprod()
Return['Buy&Hold'].iloc[0] = 1

Return['model'] = 1
for idx,val in enumerate(y_test.index) :
    if idx==0:
        Return['model'].iloc[idx] = 1
    else:
        if result['Final'].loc[val] == 0 :
            Return['model'].loc[val] = Return['model'].iloc[idx-1] +
spy['Next Return'].loc[val]
        if result['Final'].loc[val] == 1 :
            Return['model'].loc[val] = Return['model'].iloc[idx-1] +
spy['Next Return'].loc[val]
        if result['Final'].loc[val] == -1 :
            Return['model'].loc[val] = Return['model'].iloc[idx-1]

plt.figure(figsize=(20,10))
plt.plot(Return['model'],label='model')
plt.plot(Return['Buy&Hold'],label='Buy&Hold')
plt.legend()
plt.grid()
plt.show()

chart = pd.DataFrame()
chart['Price'] = spy['Adj Close'].loc[spy.index > test_date]
chart['Signal'] = result['Final'].loc[result.index > test_date]

chart['Buy'] = chart['Price'].loc[chart['Signal']==1]
chart['Sell'] = chart['Price'].loc[chart['Signal']==-1]
chart['Hold'] = chart['Price'].loc[chart['Signal']==0]

plt.figure(figsize=(20,10))
plt.plot(chart['Price'],c='gray')
plt.scatter(chart.index,chart['Hold'],c='gray',linewidths=0.5)
plt.scatter(chart.index,chart['Buy'],c='green',linewidths=0.5)
plt.scatter(chart.index,chart['Sell'],c='red',linewidths=0.5)
plt.grid()
plt.show()

```