

الگوریتم GMDH در پیش بینی قیمت سهام در بازار بورس

پروژه مقطع کارشناسی

ساجد زرین پور

استاد راهنما: فرشید مهردوست

تابستان ۱۳۹۵

بسم الله الرحمن الرحيم

چکیده

در این مقاله یک شبکه عصبی مدل GMDH و الگوریتم ژنتیک برای پیش بینی قیمت سهام ارائه شده است. برای پیش بینی قیمت سهام با شبکه عصبی مبتنی بر الگوریتم GMDH، ما از سود هر سهم (EPS)، سود پیش بینی شده هر سهم (PEPS)، سود نقدی هر سهم (DPS)، نسبت قیمت به درآمد (P/E) و نسبت درآمد به قیمت (E/P) به عنوان داده‌های ورودی و قیمت سهام به عنوان داده خروجی استفاده می‌کنیم. برای این کار، از یک پایگاه داده ساختگی استفاده کرده ایم. شبکه عصبی با ۸۵٪ داده های تجربی ساخته شده است. برای آزمودن تناسب بودن مدل، از باقی داده ها در شبکه استفاده می‌کنیم.

فهرست	
۱	مقدمه
۲	نگاهی به مفاهیم کلی سهام
۳	نگاهی به الگوریتم GMDH
۴	الگوریتم GMDH عددی
۷	شبکه عصبی GMDH
۹	پیاده سازی الگوریتم در نرم افزار متلب
۱۱	نتیجه گیری
	ضمایم:
۱۲	ضمیمه ۱. کد برنامه
۲۰	ضمیمه ۲. اطلاعات مدل
۲۶	مراجع

پیش بینی قیمت سهام یکی از وظایف اصلی سرمایه گذاران حقیقی و حقوقی می باشد. این، یک مسئله مهم در تصمیم گیری های سرمایه گذاری مالی است و در حال حاضر توجه زیادی را از سوی جامعه محققین به خود جلب کرده است. گرچه، این مسئله با توجه به این که قیمت سهام / سود آن پرنوسان و غیرثابت است، به عنوان یکی از چالش برانگیزترین مسائل قلمداد می گردد (Hall, 1994, Li et al 2003, Yaser and Atiya 1996).

تغییرات قیمت سهام یک سیستم بسیار دینامیک است که از تعدادی اصول مشتق شده است. دو روش اصلی تحلیلی، آنالیز بنیادی و آنالیز فنی اند. آنالیز بنیادی از فاکتورهای کلی اقتصادی مانند نرخ سرمایه گذاری، سرمایه مالی، نرخ های تورم، و نرخ مبادلات خارجی به عنوان وضعیت بنیادی مالی شرکت استفاده می کند. پس از تحلیل دقیق تمام این فاکتورها، تحلیل گر تصمیمی مبتنی بر خرید یا فروش یک سهم می گیرد. یک تحلیل فنی مبتنی بر داده های سری های زمانی مالی است. گرچه، سری های زمانی مالی کاملاً داده های پیچیده ای ارائه می دهند (برای مثال، روندها، تغییرات ناگهانی، فراریت دسته بندی) و این گونه سری ها اغلب غیر ثابت-اند، جایی که یک متغیر گرایش روشنی برای حرکت به سمت یک مقدار ثابت یا یک روند خطی ندارد (Cheng and Liu 2008).

هدف این مقاله کاربرد شبکه های عصبی مبتنی بر الگوریتم GMDH برای پیش بینی قیمت سهام است. برای این کار، ما از سود هر سهم (EPS)، سود پیش بینی شده هر سهم (PEPS)، سود نقدی هر سهم (DPS)، نسبت قیمت به درآمد (P/E) و نسبت درآمد به قیمت (E/P) به عنوان داده های ورودی و قیمت سهام به عنوان داده خروجی استفاده می کنیم. دیدگاه تحلیل پیشنهادی در این مقاله از شبکه عصبی مبتنی بر الگوریتم GMDH برای پیش بینی قیمت سهام استفاده می کند و می تواند کمک خوبی برای سرمایه گذاران و تحلیل گران مالی باشد. این مقاله به صورت زیر سازمان دهی شده است.

بخش ۱ به معرفی مفاهیم اقتصادی می پردازد. بخش ۲ به معرفی اجمالی الگوریتم GMDH اقدام نموده و در بخش ۳ صورت ریاضی الگوریتم GMDH را بیان می نماید. بخش ۴ به توصیف اجمالی شبکه عصبی مبتنی بر الگوریتم GMDH پرداخته و در بخش ۵ توضیحات مربوط به پیاده سازی این الگوریتم در نرم افزار متلب آورده شده است. و در انتها کد کامل برنامه در ضمیمه ۱ و اطلاعات مدل در ضمیمه ۲ آورده شده اند. کد کامل برنامه به همراه یک پایگاه داده تمثیلی برای مشاهده نحوه عملکرد برنامه در فایل های ضمیمه در اختیار علاقه مندان قرار گرفته است.

۱. نگاهی به مفاهیم کلی سهام

۱.۱ درآمد هر سهم (EPS): با محاسبه این رقم، سودی که شرکت در یک دوره مشخص به ازای یک سهم عادی به دست آورده معین می شود.

$$EPS = \frac{\text{سود پس از کسر مالیات}}{\text{تعداد سهام عادی منتشره}}$$

علت اصلی توجه به سود هر سهم (و نه سود کل شرکت) به هدف اصلی شرکت مربوط است که به حداکثر رسانیدن ثروت سهامداران می باشد.

۲.۱ درآمد پیش بینی شده هر سهم (PEPS): آخرین پیش بینی از درآمد هر سهم است.

۳.۱ سود نقدی هر سهم (DPS): به مقدار سودی که شرکت تقسیم می کند و به طور نقدی به دست سهامدار می رسد، DPS گفته می شود. اینکه چه میزان از سود تحقق یافته هر سهم به صورت نقدی توزیع شود و چه میزان در شرکت باقی بماند تصمیمی است که در مجمع عادی سالانه توسط هیئت مدیره به مجمع پیشنهاد داده می شود و سپس به تصویب سهامداران می رسد. فرمول محاسبه سود نقدی توزیع شده (DPS) در ذیل نشان داده شده است.

$$DPS = EPS - (\text{اندوخته قانونی} + \text{اندوخته طرح و توسعه} + \text{سود انباشته})$$

۴.۱ نسبت E/P و P/E سود هر سهم، E/P مبلغ سود به ازای ارزش اسمی سهام را نشان می دهد. با توجه به اینکه ارزش اسمی سهام کلیه شرکت های پذیرفته شده در بورس تهران ۱۰۰۰ ریالی بوده و ممکن است قیمت بازار انواع سهام متفاوت باشد، برای مقایسه شرکت ها با یکدیگر رقم سود نسبت به قیمت بازار تحت عنوان E/P سنجیده می شود. این نسبت بازده جاری سهام را نشان داده و در شرایط تساوی سایر عوامل، هرچه این نسبت بالاتر باشد سهام از موقعیت بهتری برخوردار خواهد بود. لازم به ذکر است هرچه میزان سود در هر سهم در مقایسه با قیمت بازار بیشتر باشد به خاطر تقاضایی که برای آن ایجاد می شود انتظار احتمال افزایش آتی قیمت بالاتر خواهد بود. خاطر نشان می شود که سهام داران در خرید سهام می توانند از نسبت قیمت به درآمد تحت عنوان P/E نیز استفاده کنند. نسبت P/E یکی از ابزارهای رایج و متداول جهت تحلیل وضعیت شرکت ها، صنعت و بازار است. این نسبت که از آن به عنوان ضریب سودآوری نیز نام برده می شود، حاصل تقسیم سهم بردرآمد هر سهم آن شرکت است و در واقع رابطه بین سهام شرکت با سود آن را نشان می دهد. همانطور که گفته شد، سود هر سهم، معیاری است که می توان با آن یک شرکت را طی چند سال متوالی از لحاظ افزایش سرمایه های شرکت مقایسه نمود ولی این متغیر، معیار مناسبی برای مقایسه چند شرکت نیست. این نقیصه در معیار P/E رفع می شود و با این معیار می توان شرکت های یک صنعت را مورد مقایسه قرار داد یا هر شرکت را با متوسط P/E صنعت مرتبط تطبیق داد. نسبت P/E بیانگر تعداد دفعاتی است که طول می کشد تا قیمت یک سهم از طریق سود آن بازیافت شود. البته تحلیل P/E براساس دوره بازگشت سرمایه باید توأم با احتیاط باشد. به دلیل این که شرکت ها در هر سال سود یکسانی ندارند و درآمد هر سهم آنها توأم با نوسان است. بنابراین در صورتی که P/E شرکتی برابر ۱۰ باشد لزوماً به این معنی نیست که ۱۰ سال طول می کشد تا سرمایه بازگشت کند. اگر شرکت در سال های بعد سود بیشتری کسب کند، دوره بازگشت سرمایه کاهش می یابد. لازم به ذکر است که برای مقایسه P/E بین شرکت ها، مقایسه در درون یک صنعت انجام شود. با فرض ثابت

بودن سایر شرایط به نظر می رسد شرکتی که در گروه خود دارای P/E کوچکتری می باشد امکان رشد قیمت سهام آن در آینده نزدیک وجود دارد.

۲. نگاهی به الگوریتم GMDH

الگوریتم پایه ای GMDH یک روش ساخت چند جمله ای با مرتبه بالا به فرم

$$y = a + \sum_{i=1}^m b_i x_i + \sum_{i=1}^m \sum_{j=1}^m c_{ij} x_i x_j + \sum_{i=1}^m \sum_{j=1}^m \sum_{k=1}^m d_{ijk} x_i x_j x_k + \sum_{i=1}^m \sum_{j=1}^m \sum_{k=1}^m \sum_{t=1}^m e_{ijkl} x_i x_j x_k x_t + \dots \quad (1.2)$$

است که m متغیر ورودی x_1, x_2, \dots, x_m را به یک متغیر خروجی y مرتبط می سازد. گرچه (۱.۲) به یک چند جمله ای با مرتبه بالای رگرسیون شباهت دارد، روشی که این چندجمله ای (که از این پس آن را چند جمله ای ایواخنکف می نامیم؛) طی آن ساخته می شود با تکنیک های متداول در آنالیز استاندارد رگرسیون تفاوت دارد. در واقع گفته می شود (اسکات و هاچینسون؛ ۱۹۷۶) که روش ساخت چندجمله ای ایواخنکف به روشی که طبیعت طی آن انتخاب طبیعی را انجام می دهد مشابه است. این روندی است که اکنون به توصیف آن خواهیم پرداخت.

۳. الگوریتم عددی GMDH

اطلاعات پایه ای که یک فرد باید برای ساخت چندجمله ای ایواخنکف جمع آوری کند یک مجموعه از n مشاهده مانند آنچه در شکل-۱ نمایش داده شده، است.

		Y						X					
مشاهدات یادگیری	y_1	x_{11}						x_{12}					
	y_2	x_{21}						x_{22}					
	\vdots	\vdots						\vdots					
	y_{nt}	$x_{nt,1}$						$x_{nt,2}$					
	\vdots	\vdots						\vdots					
مشاهدات آزمون	\vdots	\vdots						\vdots					
	\vdots	\vdots						\vdots					
	y_n	x_{n1}						x_{n2}					
		x_1						x_2					
		y						x_m					

شکل-۱. ورودیهای الگوریتم

در شکل فوق :

nt : تعداد مشاهدات در مجموعه یادگیری،

$nc=n-nt$: تعداد مشاهدات در مجموعه آزمون،

n : تعداد کل مشاهدات و

m : تعداد متغیرها.

دلیل تقسیم مشاهدات به دو مجموعه مجزا به زودی بیان می شود. حال به توصیف گام های اصلی الگوریتم GMDH می پردازیم.

گام ۱. (ساخت متغیرهای جدید $Z_1, Z_2, \dots, Z_{\binom{m}{2}}$)

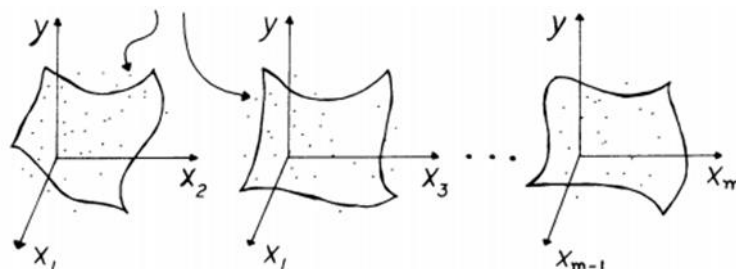
از تمام متغیرهای مستقل x_1, x_2, \dots, x_m (ستون های آرایه X) هر بار ۲ متغیر انتخاب می کنیم و برای هریک از این $\binom{m}{2}$ ترکیب، چندجمله ای تقریب کمترین مربعات به فرم

$$y = A + Bu + Cv + Du^2 + Ev^2 + Fuv \quad (۱.۳)$$

را که بهترین تقریب برای مشاهده y_i در مجموعه یادگیری می باشد را پیدا می کنیم. بدین گونه $\binom{m}{2}$ صفحه چندجمله ای مانند آنچه در شکل ۲ نشان داده شده است، می یابیم. حال،

برای هریک از این $\binom{m}{p}$ صفحه چندجمله ای نشان داده شده در شکل ۲، چند جمله ای را در n نقطه مقداربایی می کنیم.

مشاهدات مجموعه یادگیری



شکل-۲ ساختار متغیرهای جدید

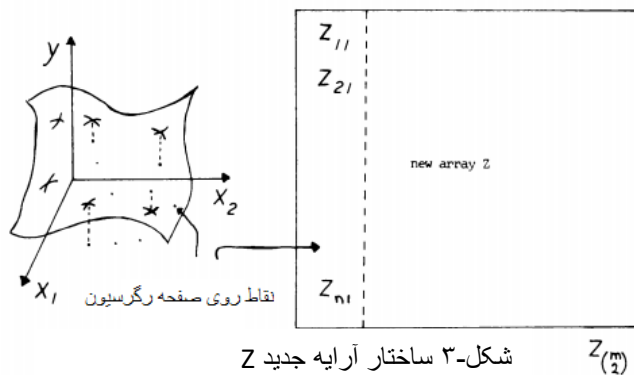
برای مثال برای اولین سطح نشان داده شده در بالا، ما چندجمله ای $y_k = A + Bx_i + Cx_j + Dx_i^2 + Ex_j^2 + Fx_ix_j$ را در نقاط $(x_{1i}, x_{1j}), (x_{2i}, x_{2j}), \dots, (x_{ni}, x_{nj})$ مقداربایی می کنیم و این n متغیر را در k -امین ستون از آرایه Z ذخیره می کنیم. $1 - \binom{m}{p}$ ستون باقی مانده به روش مشابهی ساخته می شوند (شکل ۳ را مشاهده کنید).

می توان این $\binom{m}{p}$ ستون از Z را به عنوان $\binom{m}{p}$ مشاهده از متغیرهای $Z_1, Z_2, \dots, Z_{\binom{m}{p}}$ که هریک از این متغیرهای جدید یک چندجمله ای از متغیرهای x_1, x_2, \dots, x_m اند، در نظر گرفت. به عبارت دیگر ما متغیرهای جدیدی ساختیم که برخی از آنها با متغیرهای اصلی تعویض خواهند شد. هدف نگاه داشتن آن Z_i هایی است که بهترین تخمین از متغیر خروجی y ارائه می دهند و دور ریختن اضافات است. این، جایی است که مجموعه آزمون وارد الگوریتم می شود.

گام ۲. (یافتن Z -هایی که مهم نیستند).

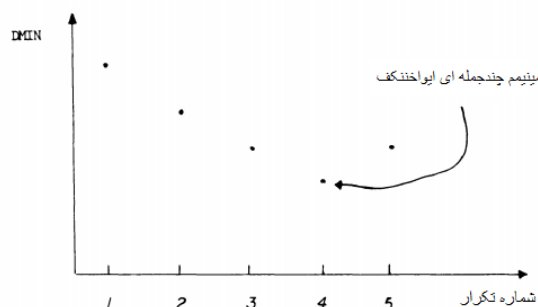
در این گام ستون هایی از X (متغیرهای قدیمی) را با ستون هایی از Z (متغیرهای جدید که چندجمله ای های مربعی از متغیرهای قدیمی اند) که بهترین پیش بینی از بردار مشاهده شده y ارائه می دهند، تعویض می کنیم. به طور دقیق تر، برای هر ستون j از Z ، ما خطای کمترین مربعات d_j را توسط $d_j^2 = \sum_{i=nt+1}^n (y_i - z_{ij})^2$ محاسبه می کنیم، سپس ستون های Z را براساس افزایش کمترین خطای مربعات، d_j ، مرتب می کنیم و پس از آن ستون هایی را که به ازای آنها $d_j < M$ (یک عدد از پیش تعیین شده است)، برقرار است را انتخاب می کنیم تا با ستون های X تعویض شوند. (این ستون ها از Z ، آرایه X جدید را می سازند).

در این جا ذکر یک نکته ضروری است. تعداد این ستون های تعویض شده (مثلا m_i) ممکن است از تعداد ستون های اصلی (m) بیشتر یا کمتر باشد. همچنین توجه می کنیم که خطای d_j روش کمترین مربعات، برابر مجموع خطاهای مجموعه مشاهدات آزمون می باشد.



گام ۳. (آزمون همگرایی)

از گام دوم، کوچکترین d_j را انتخاب می کنیم و آن را DMIN می نامیم و نمودار آن را مانند آنچه در شکل ۴ آمده، رسم می کنیم. اگر مقدار DMIN از مقدار DMIN قبلی کوچکتر باشد (اولین بار که وارد این گام می شویم درستی این مطلب را فرض می کنیم)،



شکل-۴ معیار توقف

باز می گردیم و گام های اول و دوم را تکرار می کنیم. اما اگر مقدار جدید DMIN از مقدار قبلی آن بیشتر بود روند را متوقف می کنیم و از نتایج مربوط به DMIN قبلی استفاده می کنیم.

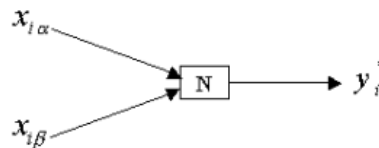
رگرسیون استاندارد تنها از میانگین خطای مربعات به عنوان معیار استفاده می کند و بنابراین تخمین اینکه مدل بسیار ساده یا بسیار پیچیده است غیرممکن می باشد. با محک زدن چندجمله ای ایواخنکف با مجموعه آزمون، بدست آوردن یک چندجمله ای یکتا با پیچیدگی بهینه ممکن است. گرچه تئوری الگوریتم GMDH به ما اجازه اثبات یکتایی مقدار مینیمم DMIN را نمی دهد، اما تمامی نتایج تجربی درستی این موضوع را تایید می کند. ذات منحنی DMIN اینگونه است که ابتدا به یک مقدار مینیمم کاهش می یابد و پس از آن شروع به افزایش می کند. برای مثال در شکل ۴ ما روند را پس از تکرار پنجم متوقف می کنیم. این، پایان الگوریتم ایواخنکف است. حال اولین ستون Z شامل مقادیر \bar{y}_i از چندجمله ای ایواخنکف

$$\begin{aligned}\bar{y}_1 &= a + \sum_{i=1}^m b_i x_{1i} + \sum_{i=1}^m \sum_{j=1}^m c_{ij} x_{1i} x_{1j} + \dots \\ &\vdots \\ \bar{y}_n &= a + \sum_{i=1}^m b_i x_{ni} + \sum_{i=1}^m \sum_{j=1}^m c_{ij} x_{ni} x_{nj} + \dots\end{aligned}$$

مقداریابی شده در n نقطه داده اصلی می باشد. به عبارت دیگر ستون اول باید به میزان زیادی به متغیر وابسته y همگرا باشد. برای یافتن مقادیر مطلوب ضرایب a, b_i, c_{ij}, \dots در چندجمله ای ایواخنکف (۱.۱) باید تمامی سطوح رگرسیون به فرم (۱.۳) را که در هر تکرار محاسبه می شوند، ذخیره شوند و به طور منظم، این چندجمله ای های متشکل از چندجمله ای های متشکل از چندجمله ای ها مقداریابی شوند، و این روند تا زمان رسیدن به چندجمله ای ایواخنکف ادامه داده شود. این، یک مسئله غیر بدیهی برنامه نویسی است اما بهر حال قابل حل است.

۴. شبکه عصبی GMDH

شبکه عصبی GMDH، شبکه ای خود سازمان ده و یک سویه است که از چندین لایه و هر لایه نیز از چندین نرون تشکیل شده است. تمامی نرون ها از ساختار مشابهی برخوردارند؛ به طوری که دارای دو ورودی و یک خروجی هستند. هر نرون با پنج وزن و یک بایاس عمل پردازش را میان داده های ورودی و خروجی براساس رابطه زیر برقرار می کند.



$$y_{ik}^* = N(x_{i\alpha}, x_{i\beta}) = b^k + w_1^k x_{i\alpha} + w_2^k x_{i\beta} + w_3^k x_{i\alpha}^2 + w_4^k x_{i\beta}^2 + w_5^k x_{i\alpha} x_{i\beta} \quad (1.4)$$

که در آن

$i = 1, 2, \dots, n$ (که n تعداد داده ورودی و خروجی است)،

$k = 1, 2, \dots, \binom{m}{2}$ و

$\alpha, \beta \in \{1, 2, 3, \dots, m\}$ (که m تعداد نرون های لایه قبلی است).

وزن ها بر اساس روش خطای کمترین مربعات محاسبه شده و سپس به منزله مقادیر مشخص و ثابت، در داخل هر نرون جایگذاری می شود. ویژگی ممتاز این نوع شبکه آن است که نرون های مرحله قبلی یا لایه قبلی، عامل و مولد تولید نرون های جدید به تعداد $\binom{m}{p} = \frac{m(m-1)}{p}$ هستند و از میان نرون های تولید شده، به اجبار تعدادی از آنها حذف شده تا بدین وسیله از واگرایی شبکه جلوگیری شود.

نرون هایی که برای ادامه و گسترش شبکه باقی می ماندند، ممکن است برای ایجاد فرم همگرایی شبکه و عدم ارتباط آنها با نرون لایه آخر حذف شوند. به این لایه ها در اصطلاح نرون غیر فعال می گویند. معیار گزینش و حذف مجموعه ای از نرون ها در یک لایه، درصد مجموع مربعات خطا (d_j^2) میان مقادیر خروجی واقعی (y_j) و خروجی نرون j -ام (y_{ij}^*) است.

$$d_j^2 = \sum_{i=nt+1}^n (y_i - y_{ij}^*)^2 \quad (2.4)$$

که در آن $j \in \{1, 2, 3, \dots, \binom{m}{p}\}$ و m تعداد نرون های گزینش شده در لایه قبلی است.

نگاشتی که بین متغیرهای ورودی و خروجی توسط این نوع از شبکه های عصبی برقرار می شود، به صورت تابع غیرخطی ولترا و به شکل رابطه شماره (۱.۳) است. ساختاری که برای نرون ها در نظر گرفته شده، به صورت فرم خلاصه شده دو متغیره درجه دوم زیر است.

$$y_i = f(x_{ip}, x_{iq}) = a. + a_1 x_{ip} + a_2 x_{iq} + a_3 x_{ip} x_{iq} + a_4 x_{ip}^2 + a_5 x_{iq}^2 \quad (3.4)$$

تابع f دارای شش ضریب مجهول است که به ازای تمام نمونه های دو متغیر وابسته به سیستم $\{(x_{ip}, x_{iq}), i = 1, 2, \dots, n\}$ خروجی مطلوب $\{(y_i), i = 1, 2, \dots, n\}$ را برآورد می کند. تابع f را براساس قاعده خطای کمترین مربعات پایه ریزی می کنیم.

$$\min \sum_{k=1}^n [(f(x_{ki}, x_{kj}) - y_i)^2] \quad (4.4)$$

براین اساس دستگاه معادله ای را حل می کنیم که دارای شش مجهول و n معادله است.

$$\begin{cases} a. + a_1 x_{1p} + a_2 x_{1q} + a_3 x_{1p} x_{1q} + a_4 x_{1p}^2 + a_5 x_{1q}^2 = y_1 \\ a. + a_1 x_{2p} + a_2 x_{2q} + a_3 x_{2p} x_{2q} + a_4 x_{2p}^2 + a_5 x_{2q}^2 = y_2 \\ \vdots \\ a. + a_1 x_{np} + a_2 x_{nq} + a_3 x_{np} x_{nq} + a_4 x_{np}^2 + a_5 x_{nq}^2 = y_n \end{cases} \quad (5.4)$$

دستگاه معادله فوق را می توان به شکل ماتریسی $Aa=Y$ نمایش داد که در آن

$$A = \begin{bmatrix} 1 & x_{1p} & x_{1q} & x_{1p}x_{1q} & x_{1p}^2 & x_{1q}^2 \\ 1 & x_{2p} & x_{2q} & x_{2p}x_{2q} & x_{2p}^2 & x_{2q}^2 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_{np} & x_{nq} & x_{np}x_{nq} & x_{np}^2 & x_{nq}^2 \end{bmatrix}$$

$$Y = [y_1 \ y_2 \ y_3 \ \dots \ y_n] \text{ و } a = [a_1 \ a_2 \ a_3 \ a_4 \ a_5 \ a_6 \ a_7]^T ,$$

برای تعیین دقت مدل تعدادی از معیارهای آماری استفاده شده است. این معیارها عبارتند از: جذر میانگین مربعات خطا (RMSE)، و قدرمطلق انحراف معیار (MAD).

$$RMSE = \left[\frac{\sum_{i=1}^M (Y_{i(model)} - Y_{i(actual)})^2}{M} \right]^{1/2} ,$$

$$MAD = \frac{\sum_{i=1}^M |Y_{i(model)} - Y_{i(actual)}|}{M} .$$

۵. پیاده سازی الگوریتم در نرم افزار متلب

برنامه از نقطه شروع استاندارد (تابع main) آغاز و هدایت می گردد. پس از بارگذاری پایگاه داده (در مثال پیاده سازی شده فایل Indice.mat) مجموعه های یادگیری و ناظر مشخص می شوند. در این مثال ۸۵٪ داده ها به مجموعه یادگیری (ptrain=0.85) و باقی به مجموعه آزمون داده شده است. پس از آن شبکه عصبی به کمک تابع GMDH و با استفاده از مجموعه یادگیری ساخته می شود، و در آخر مقداربایی ها صورت می گیرند و نتایج به صورت نمودارهایی نمایش داده می شوند.

بارگذاری پایگاه داده و مقداردهی اولیه (در پایگاه داده، آرایه Y' در جدول Targets و آرایه X' در جدول Inputs ذخیره شده اند. سطرها به ترتیب حاوی مقادیر متغیرهای EPS,PEPS,DPS,E/P,P/E اند).

```
data = load('Indice.mat');
Inputs = data.Inputs;
Targets = data.Targets;
nData = size(Inputs,2);
Perm = randperm(nData);
```

مشخص کردن مجموعه یادگیری

```
pTrain = 0.85;
nTrainData = round(pTrain*nData);
TrainInd = Perm(1:nTrainData);
TrainInputs = Inputs(:,TrainInd);
```

```
TrainTargets = Targets(:,TrainInd);
```

مشخص کردن مجموعه آزمون

```
pTest = 1 - pTrain;
```

```
nTestData = nData - nTrainData;
```

```
TestInd = Perm(nTrainData+1:end);
```

```
TestInputs = Inputs(:,TestInd);
```

```
TestTargets = Targets(:,TestInd);
```

ساخت شبکه عصبی با استفاده از مجموعه یادگیری

```
params.MaxLayerNeurons = 15; % Maximum Number of Neurons in a Layer
```

```
params.MaxLayers = 4; % Maximum Number of Layers
```

```
params.alpha = 0.6; % Selection Pressure (in Layers)
```

```
params.pTrain = 0.85; % Train Ratio (vs. Validation Ratio)
```

```
gmdh = GMDH(params, TrainInputs, TrainTargets);
```

ارزشیابی شبکه GMDH

```
Outputs = ApplyGMDH(gmdh, Inputs);
```

```
TrainOutputs = Outputs(:,TrainInd);
```

```
TestOutputs = Outputs(:,TestInd);
```

نمایش نتایج

```
disp('Type "gmdh.Layers" to see the layers" info.');
```

```
disp(' ');
```

```
figure;
```

```
PlotResults(TrainTargets, TrainOutputs, 'Train Data');
```

```

figure;
PlotResults(TestTargets, TestOutputs, 'Test Data');

figure;
PlotResults(Targets, Outputs, 'All Data');

if ~isempty(which('plotregression'))
    figure;
    plotregression(TrainTargets, TrainOutputs, 'Train Data', ...
        TestTargets, TestOutputs, 'TestData', ...
        Targets, Outputs, 'All Data');
end

```

برای مشاهده کد کامل برنامه به انتهای مقاله مراجعه فرمایید. علاوه براین تمام کدها در فایل های ضمیمه موجود اند.

۶. نتیجه گیری

در این مقاله ما رابطه بین سود سهام و قیمت آن را بررسی کردیم. مدل کردن نرم افزاری پیش بینی قیمت سهام برای تمامی سرمایه گذاران و تحلیل گران برای کاهش ریسک سرمایه گذاری و افزایش منافع سهام داران مفید است.

```
%  
% Copyright (c) 2015, Yarpiz (www.yarpiz.com)  
% All rights reserved. Please read the "license.txt" for license  
terms.  
%  
% Project Code: YPML113  
% Project Title: Implementation of Group Method of Data Handling  
in MATLAB  
% Publisher: Yarpiz (www.yarpiz.com)  
%  
% Developer: S. Mostapha Kalami Heris (Member of Yarpiz Team)  
%  
% Contact Info: sm.kalami@gmail.com, info@yarpiz.com  
%  
  
clc;  
clear;  
close all;  
  
%% Load Data  
  
data = load('Indice.mat');  
Inputs = data.Inputs;  
Targets = data.Targets;  
  
nData = size(Inputs,2);  
Perm = randperm(nData);  
  
% Train Data  
pTrain = 0.85;  
nTrainData = round(pTrain*nData);  
TrainInd = Perm(1:nTrainData);  
TrainInputs = Inputs(:,TrainInd);  
TrainTargets = Targets(:,TrainInd);  
  
% Test Data  
pTest = 1 - pTrain;  
  
  
nTestData = nData - nTrainData;
```



```

TestInd = Perm(nTrainData+1:end);
TestInputs = Inputs(:,TestInd);
TestTargets = Targets(:,TestInd);

%% Create and Train GMDH Network

params.MaxLayerNeurons = 15;    % Maximum Number of Neurons in a
Layer
params.MaxLayers = 4;           % Maximum Number of Layers
params.alpha = 0.6;             % Selection Pressure (in Layers)
params.pTrain = 0.85;           % Train Ratio (vs. Validation Ratio)
gmdh = GMDH(params, TrainInputs, TrainTargets);

%% Evaluate GMDH Network

Outputs = ApplyGMDH(gmdh, Inputs);
TrainOutputs = Outputs(:,TrainInd);
TestOutputs = Outputs(:,TestInd);

%% Show Results

disp('Type "gmdh.Layers" to see the layers'' info. ');
disp(' ');

figure;
PlotResults(TrainTargets, TrainOutputs, 'Train Data');

figure;
PlotResults(TestTargets, TestOutputs, 'Test Data');

figure;
PlotResults(Targets, Outputs, 'All Data');

if ~isempty(which('plotregression'))
    figure;
    plotregression(TrainTargets, TrainOutputs, 'Train Data', ...
                    TestTargets, TestOutputs, 'TestData', ...
                    Targets, Outputs, 'All Data');
end

```

```
function gmdh = GMDH(params, X, Y)

    disp('Training GMDH:');

    MaxLayerNeurons = params.MaxLayerNeurons;
    MaxLayers = params.MaxLayers;
    alpha = params.alpha;

    nData = size(X,2);

    % Shuffle Data
    Permutation = randperm(nData);
    X = X(:,Permutation);
    Y = Y(:,Permutation);

    % Divide Data
    pTrainData = params.pTrain;
    nTrainData = round(pTrainData*nData);
    X1 = X(:,1:nTrainData);
    Y1 = Y(:,1:nTrainData);
    pTestData = 1-pTrainData;
    nTestData = nData - nTrainData;
    X2 = X(:,nTrainData+1:end);
    Y2 = Y(:,nTrainData+1:end);

    Layers = cell(MaxLayers, 1);

    Z1 = X1;
    Z2 = X2;

    for l = 1:MaxLayers

        L = GetPolynomialLayer(Z1, Y1, Z2, Y2);

        if l>1
            if L(1).RMSE2 > Layers{l-1}(1).RMSE2
                break;
            end
        end

        ec = alpha*L(1).RMSE2 + (1-alpha)*L(end).RMSE2;
        ec = max(ec, L(1).RMSE2);
        L = L([L.RMSE2] <= ec);
    end
end
```

```

    if numel(L) > MaxLayerNeurons
        L = L(1:MaxLayerNeurons);
    end

    if l==MaxLayers && numel(L)>1
        L = L(1);
    end

    Layers{l} = L;

    Z1 = reshape([L.Y1hat],nTrainData,[])';
    Z2 = reshape([L.Y2hat],nTestData,[])';

    disp(['Layer ' num2str(l) ': Neurons = ' num2str(numel(L))
', Min Error = ' num2str(L(1).RMSE2)]);

    if numel(L)==1
        break;
    end

end

Layers = Layers(1:l);

gmdh.Layers = Layers;

disp(' ');

end

```

```
function Yhat = ApplyGMDH(gmdh, X)

    nLayer = numel(gmdh.Layers);

    Z = X;
    for l=1:nLayer
        Z = GetLayerOutput(gmdh.Layers{l}, Z);
    end
    Yhat = Z;

end

function Z = GetLayerOutput(L, X)

    m = size(X,2);
    N = numel(L);
    Z = zeros(N,m);

    for k=1:N
        vars = L(k).vars;
        x = X(vars,:);
        Z(k,:) = L(k).f(x);
    end

end
```

```
function PlotResults(Targets, Outputs, Title)

    Errors = Targets - Outputs;
    MSE = mean(Errors.^2);
    RMSE = sqrt(MSE);
    ErrorMean = mean(Errors);
    ErrorStd = std(Errors);

    subplot(2,2,[1 2]);
    plot(Targets, 'r');
    hold on;
    plot(Outputs);
    legend('Targets', 'Outputs');
    ylabel('Targets and Outputs');
    grid on;
    title(Title);

    subplot(2,2,3);
    plot(Errors);
    title(['MSE = ' num2str(MSE) ', RMSE = ' num2str(RMSE)]);
    ylabel('Errors');
    grid on;

    subplot(2,2,4);
    histfit(Errors, 50);
    title(['Error Mean = ' num2str(ErrorMean) ', Error StD = '
num2str(ErrorStd)]);

end
```

```
function p = FitPolynomial(x1, Y1, x2, Y2, vars)
    X1 = CreateRegressorsMatrix(x1);
    c = Y1*pinv(X1);

    Y1hat = c*X1;
    e1 = Y1- Y1hat;
    MSE1 = mean(e1.^2);
    RMSE1 = sqrt(MSE1);
    MAD1=mad(abs(e1));%%

    f = @(x) c*CreateRegressorsMatrix(x);

    Y2hat = f(x2);
    e2 = Y2- Y2hat;
    MSE2 = mean(e2.^2);
    RMSE2 = sqrt(MSE2);
    MAD2=mad(abs(e2));%%

    p.vars = vars;
    p.c = c;
    p.f = f;
    p.Y1hat = Y1hat;
    p.MSE1 = MSE1;
    p.RMSE1 = RMSE1;
    p.Y2hat = Y2hat;
    p.MSE2 = MSE2;
    p.RMSE2 = RMSE2;
    %-----
    p.MAD1=MAD1;
    p.MAD1=MAD1;

end

function X = CreateRegressorsMatrix(x)

    X = [ones(1,size(x,2))
        x(1,:)
        x(2,:)
        x(1,:).^2
        x(2,:).^2
        x(1,:).*x(2,:)];

end
```

```
function L = GetPolynomialLayer(X1, Y1, X2, Y2)

    n = size(X1,1);

    N = n*(n-1)/2;

    template =
    FitPolynomial(rand(2,3),rand(1,3),rand(2,3),rand(1,3),[]);

    L = repmat(template, N, 1);

    k = 0;
    for i=1:n-1
        for j=i+1:n
            k = k+1;
            L(k) = FitPolynomial(X1([i j],:), Y1, X2([i j],:), Y2,
[i j]);
        end
    end

    [~, SortOrder] = sort([L.RMSE2]);

    L = L(SortOrder);

end
```

ضمیمه ۲. اطلاعات مدل

۱.۲ اطلاعات آماری و جداول

مدل ما از سه لایه تشکیل شده است. اطلاعات هر کدام از این لایه ها به قرار زیر است:

لایه اول

$$y_1 = 593.1933 + 10.6172 x_1 - 10.9118 x_2 - 0.2618 x_1 x_2 - 0.2413 x_1^2 + 0.5034 x_2^2$$

$$y_2 = 623.4985 - 0.3098 x_1 - 10.607 x_2 + 0.0002 x_1 x_2 + 0.0079 x_1^2 - 0.0001 x_2^2$$

$$y_3 = 622.6444 - 0.3045 x_1 - 10.572 x_2 + 0.0002 x_1 x_2 + 0.0078 x_1^2 - 0.0001 x_2^2$$

$$y_4 = 592.4597 - 0.2418 x_1 - 10.727 x_2 + 0.0002 x_1 x_2 + 0.0079 x_1^2 - 0.0001 x_2^2$$

نتیجه اعمال بر مجموعه داده های یادگیری

تابع	Mad	RMSE
y_1	۲۲,۹۶۲۲	۲۶,۹۱۱۱
y_2	۲۲,۶۴۵۳	۲۹,۲۴۲۱
y_3	۲۲,۴۷۸۶	۲۹,۴۵۸۳
y_4	۲۲,۱۹۷۴	۲۹,۷۸۷۲

نتیجه اعمال بر مجموعه داده های ناظر

تابع	Mad	RMSE
y_1	۲۲,۰۵۹۸	۳۰,۷۰۶۵
y_2	۲۰,۸۳۳۴	۳۲,۰۰۳۷
y_3	۲۰,۶۲۷۸	۳۲,۰۸۲۲
y_4	۲۰,۵۰۳۶	۳۲,۳۴۳۷

لایه دوم

$$y_1 = [-5.2240 + 0.0094 x_1 + 0.112 x_2 + 0.0000 x_1 x_2 + 0.0000 x_1^2 - 0.0001 x_2^2] * 10^3$$

$$y_1 = [-5.2866 + 0.0098 x_1 + 0.109 x_2 + 0.0000 x_1 x_2 + 0.0000 x_1^2 - 0.0001 x_2^2] * 10^3$$

$$y_1 = [-5.0573 + 0.0102 x_1 + 0.0097 x_2 + 0.0000 x_1 x_2 + 0.0000 x_1^2 - 0.0001 x_2^2] * 10^3$$

نتیجه اعمال بر مجموعه داده های یادگیری

تابع	Mad	RMSE
y_1	۳۰,۴۰۰۸	۱۸,۵۲۸۵
y_2	۳۰,۳۸۷۸	۱۸,۵۳۳۸
y_3	۳۰,۰۶۴۸	۱۹,۹۷۰۵

نتیجه اعمال بر مجموعه داده های ناظر

تابع	Mad	RMSE
y_1	۲۶,۲۸۳۴	۲۱,۲۷۳۱
y_2	۲۶,۳۸۶۳	۲۱,۲۸۵۱
y_3	۲۶,۶۷۸۷	۲۱,۳۹۴۹

لایه سوم

$$y_1 = [-1.1547 + 0.1467 x_1 - 0.1413 x_2 - 0.0102 x_1 x_2 - 0.0099 x_1^2 + 0.0201 x_2^2] * 10^3$$

نتیجه اعمال بر مجموعه داده های یادگیری

تابع	Mad	RMSE
y_1	۳۱,۲۵۲۳	۱۷,۰۱۱۵

نتیجه اعمال بر مجموعه داده های ناظر

تابع	Mad	RMSE
y_1	۲۷,۰۹۰۹	۱۹,۶۰۲۷

و در نهایت چند جمله ای ایواخنکو عبارت است از

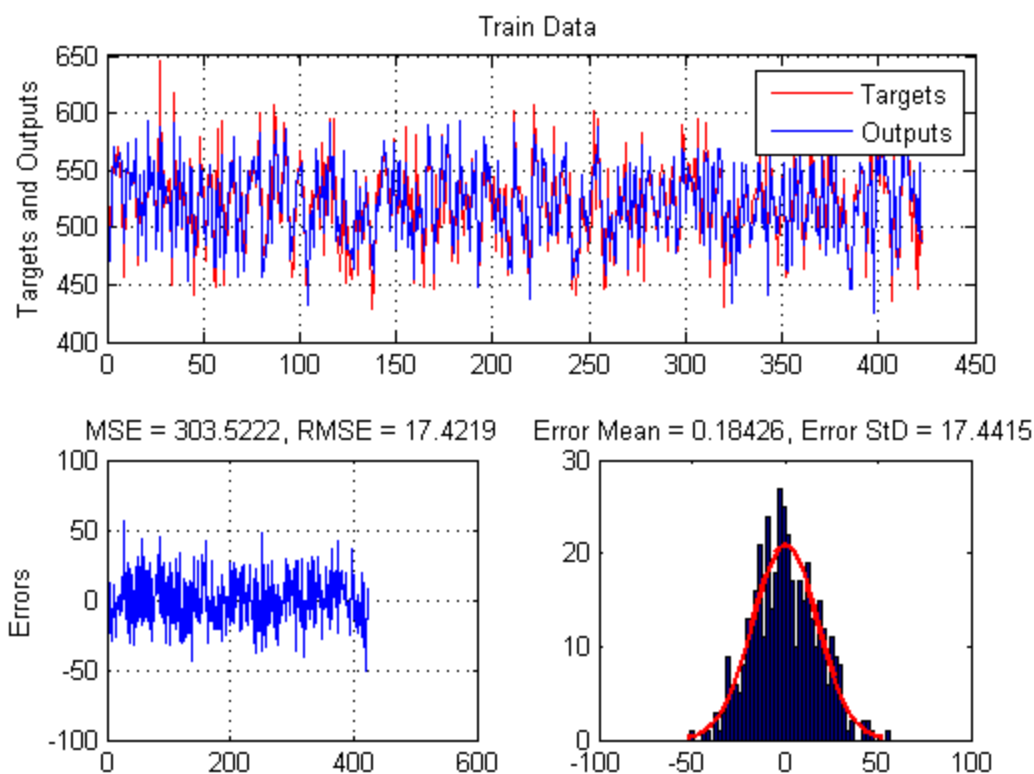
$$y = 0.0019 + 0.9231 x_1 - 0.9231 x_2 + 607.6460 x_1 x_2 + 607.5511 x_1^2 + 607.7559 x_2^2$$

که متغیرها براساس لایه آخر شبکه و ضرایب از حل دستگاه (۳.۵) توسط قطعه کد زیر بدست آمده اند.

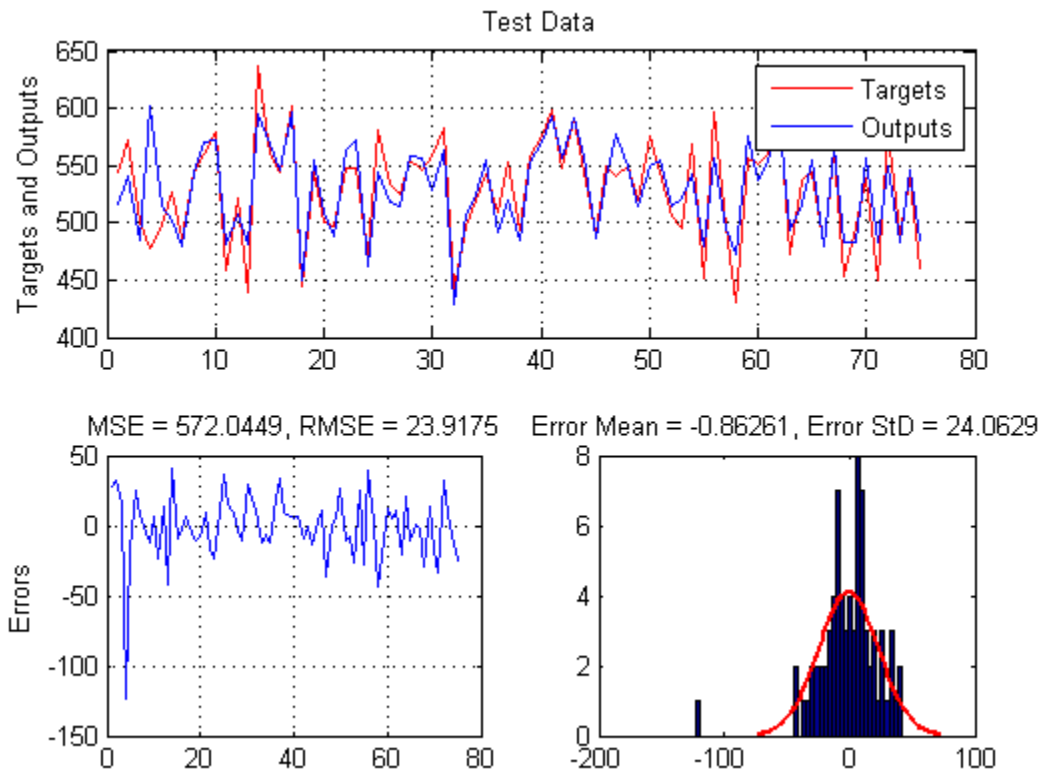
```
A=[ones(1,498)',Inputs(1,:) ',Inputs(2,:) ',Inputs(1,:) '.*Inputs(2,:) ',Inputs(1,
,:) '.*^2,Inputs(2,:) '.*^2]
disp('coefficients :')
a=Outputs'\A
```

RMSE	Mad	تابع
۲۰,۱۴۴۴	۹,۵۳۹۹	y

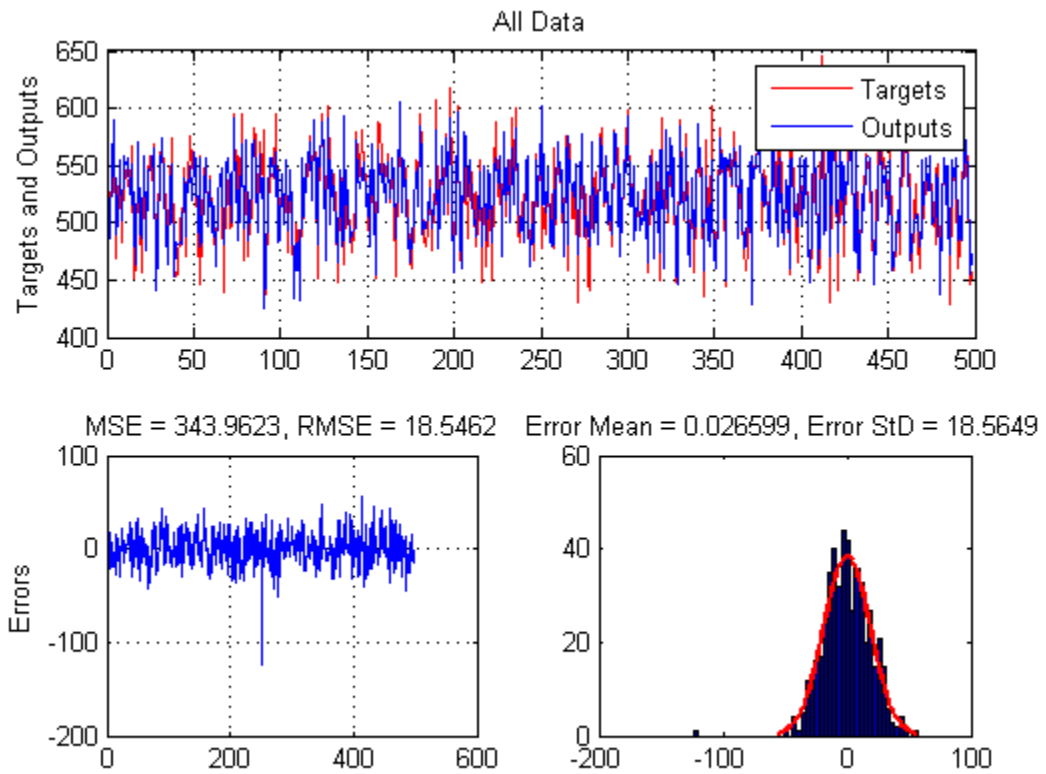
۲.۲ نمودارهای مدل



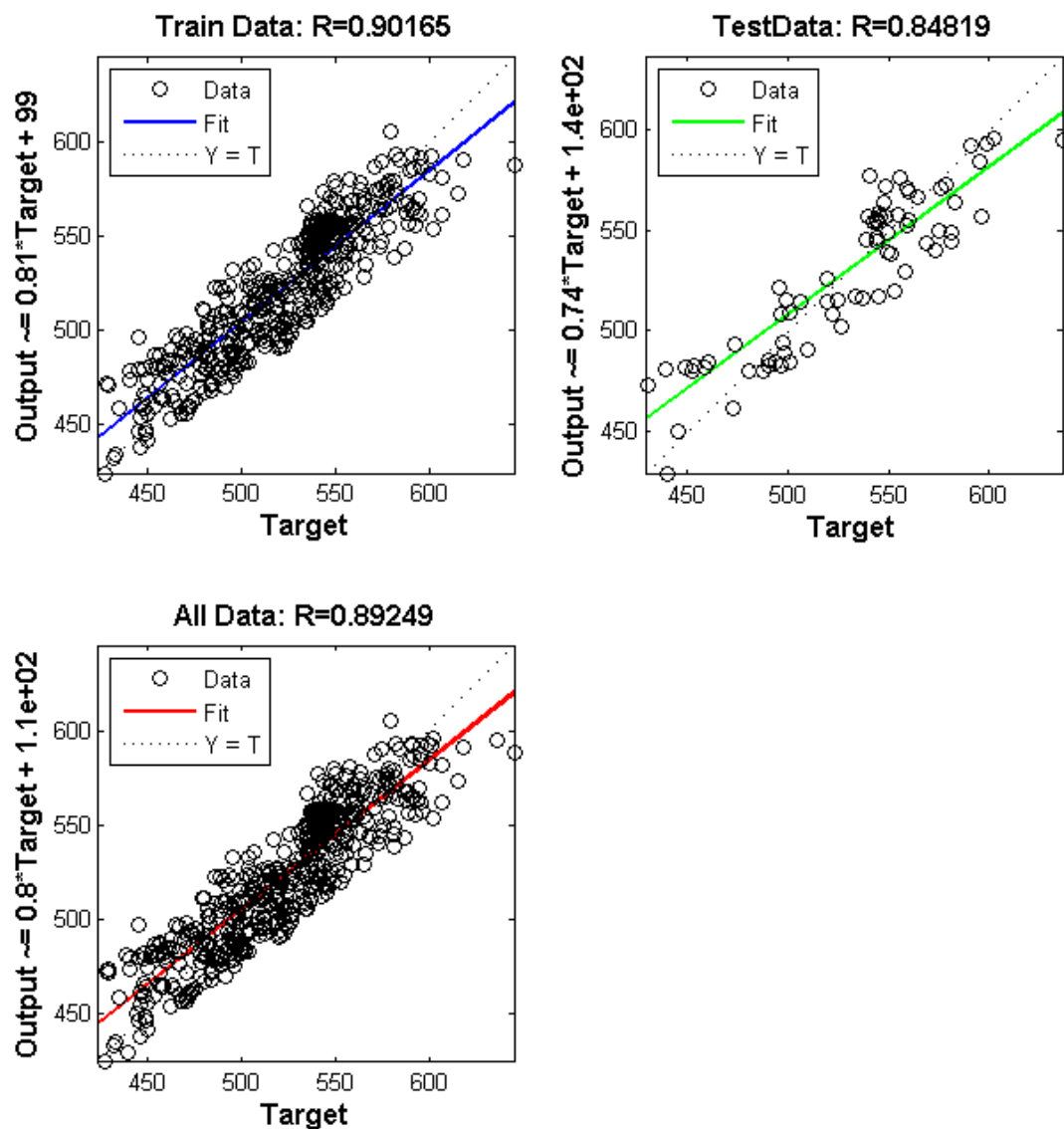
نمودار ۱.۲.۲ (نمودار بالا) مقادیر واقعی (رنگ قرمز) و مدلسازی شده (رنگ آبی) قیمت سهام به ازای مجموعه اندیس، (نمودار پایین-چپ) خطای مدل به ازای مجموعه اندیس، (نمودار پایین-راست) توزیع خطا؛ برای مجموعه یادگیری.



نمودار ۲.۲.۲ (نمودار بالا) مقادیر واقعی (رنگ قرمز) و مدلسازی شده (رنگ آبی) قیمت سهام به ازای مجموعه اندیس، (نمودار پایین-چپ) خطای مدل به ازای مجموعه اندیس، (نمودار پایین-راست) توزیع خطا؛ برای مجموعه ناظر.



نمودار ۳.۲.۲ (نمودار بالا) مقادیر واقعی (رنگ قرمز) و مدلسازی شده (رنگ آبی) قیمت سهام به ازای مجموعه اندیس، (نمودار پایین-چپ) خطای مدل به ازای مجموعه اندیس، (نمودار پایین-راست) توزیع خطا؛ برای کل داده ها.



نمودار ۴.۲.۲: نمودار رگرسیون مقادیر خروجی مدل به ازای مجموعه های یادگیری (بالا-چپ)، ناظر (بالا-راست) و کل داده ها (پایین).

مراجع :

۱. The GMDH Algorithm of Ivakjnenko, Stanle J.Farlow, The American Statistician, Vol.35, No.4(Nov.,1981),210-215

۲. Applying GMDH-Type Nural Network and Genetic Algorithm for Stock Price Prediction Of Iranian Cement Sector, Saeed Fallahi, Meysam Shaverdi, Vahab Bashiri, Applications and Applied Mathematics, Vol.6, Issue 2(December 2011), pp.572-591, ISSN: 1932-9466

۳. الگوسازی و پیش بینی EPS شرکت های پذیرفته شده در بورس اوراق بهادار تهران با رویکرد شبکه عصبی GMDH، علی اصغر انواری رستمی، عادل آذر، محمد نوروزی. مجله بررسی های حسابداری و حسابرسی، دانشکده مدیریت دانشگاه تهران، دوره ۲۰، شماره ۱، بهار ۱۳۹۲، صص ۱۸-۱.