

مقاله علمی - پژوهشی

بررسی تأثیر پیش‌پردازش داده‌ها و عوامل مدل‌سازی برنامه‌ریزی بیان ژن در دقت پیش‌بینی سری‌های زمانی

مریم صالحی^۱، سید احسان فاطمی^{۲*}

تاریخ دریافت: ۱۳۹۹/۱۰/۳۰ تاریخ پذیرش: ۱۴۰۰/۰۱/۲۴

چکیده

سری زمانی هیدرولوژیک عاملی وابسته به زمان است که یافتن نحوه تغییرات و پیش‌بینی آن مهم‌ترین هدف تجزیه و تحلیل سری‌های زمانی است. هدف این تحقیق بررسی هم‌زمان خصوصیات سری زمانی و پیش‌پردازش آن‌ها و پارامترهای مهم مدل برنامه‌ریزی بیان ژن جهت پیش‌بینی‌های با دقت بالا در مراحل آموزش و صحت‌سنجی است. در این پژوهش از سری‌های زمانی عمق آب زیرزمینی ایستگاه دشت چمچمال واقع در استان کرمانشاه با دوره زمانی ۱۲ ساله و اقلیم کوهستانی و سری زمانی ماهانه دمای آلاسکا با دوره زمانی ۵۰ ساله و اقلیم سرد و خشک استفاده شده است. برای مدل‌سازی سری‌های زمانی مذکور از نرم‌افزار Genexprotools 5.0 استفاده شده است. نتایج این تحقیق، نشان داد تناوبی بودن خصوصیات داده موجود در سری زمانی دما، سبب بروز نتایج همبستگی بالای ۹۰٪ در مراحل مختلف آموزش و صحت‌سنجی گردید به طوری که اثر پارامترهای مختلف بیان ژن کمتر از ۱۰ درصد در بهبود نتایج است. از سوی دیگر با بررسی سری زمانی عمق آب زیرزمینی که فاقد خصوصیت تناوبی و دارای شکل ACF نزولی است، نتایج پیش‌بینی مدل GEP با هر پارامتر تأثیرگذار بیان ژن، R بیش از ۴۴٪ در مرحله صحت‌سنجی حاصل نشد. این بدان معنی است که پیش‌پردازش سری زمانی اثرگذاری بیشتری در نتایج پیش‌بینی دارد. به طوری که با حذف ترم تناوب نتایج پیش‌بینی در همه مراحل مدل‌سازی به طرز معنی‌داری کاهش می‌یابد. در این حالت بهترین R برای قسمت صحت‌سنجی ۵۰ درصد است.

واژه‌های کلیدی: بیان ژن، پیش‌بینی، پیش‌پردازش، تناوب، سری زمانی

مقدمه

هستند توانایی الگوبندی فرآیندهای کاملاً غیرخطی و پویا را دارند (زمانی و همکاران، ۱۳۹۳). روش برنامه‌ریزی بیان ژن، در سال ۱۹۹۹ توسط فریرا ابداع شد. این روش ترکیبی از روش‌های GP^۳ و GA^۴ بوده که در آن، کروموزوم‌های خطی و ساده با طول ثابت، مشابه با آنچه در الگوریتم ژنتیک استفاده می‌شود و ساختارهای شاخه‌ای با اندازه‌ها و اشکال متفاوت، مشابه با درختان تجزیه در برنامه‌ریزی ژنتیک، ترکیب می‌شوند و از آنجایی که در این روش تمام ساختارهای شاخه‌ای با اندازه‌ی اشکال متفاوت، در کروموزوم‌های خطی با طول ثابت کدگذاری می‌شوند، معادل این است که در این روش فنوتیپ و ژنوتیپ از هم جدا می‌شوند و سیستم می‌تواند از تمام مزایای تکاملی به سبب وجود آن‌ها بهره‌مند شود. اکنون با وجود اینکه فنوتیپ

داده‌های هواشناسی و هیدرولوژیکی به‌عنوان اطلاعات پایه و اساسی در طراحی و مدیریت پروژه‌های منابع آب به کار می‌رود. بررسی روند و ایستایی در سری‌های زمانی هیدرولوژیکی می‌تواند در تفسیر رابطه بین فرآیندهای هیدرولوژیکی و تغییرات محیطی در مناطق مورد مطالعه کمک مؤثری داشته باشد (Salas et al., 1996). برنامه‌ریزی ژنتیک یک روش پرکاربرد در منابع آب و هیدرولوژی است. الگوریتم‌های تکاملی که برنامه‌ریزی ژنتیک نیز عضوی از آن‌ها

۱- کارشناس ارشد مهندسی منابع آب، پردیس کشاورزی و منابع طبیعی، دانشگاه رازی، کرمانشاه، ایران
۲- استادیار گروه مهندسی آب، پردیس کشاورزی و منابع طبیعی، دانشگاه رازی، کرمانشاه، ایران

* - نویسنده مسئول: (Email: e_fatemi78@yahoo.com)

DOR: 20.1001.1.20087942.1400.15.3.9.3

3- Genetic programming

4- Genetic algorithm

توانایی این روش در پیش‌بینی بارش رواناب حوضه مذکور است. پیش‌بینی نقطه شبنم روزانه با استفاده از برنامه‌ریزی بیان ژن (GEP) نیز نشان داد که با داده‌های ثبت‌شده در دو ایستگاه رودخانه ليقوان چای در آذربایجان شرقی و ایستگاه آب‌سنجی و نیار در آبخیز آجی چای و روش GEP، به‌خوبی می‌توان نقطه شبنم روزانه را پیش‌بینی کرد (مهدی زاده، ۱۳۹۶).

ادیب و همکاران جریان رودخانه کسلیان را با استفاده برنامه‌ریزی بیان ژن پیش‌بینی کردند. نتایج نشان داد که روش‌های مختلف نرمال‌سازی و حذف روند داده‌ها و حذف خصوصیات تناوبی داده‌ها به‌منظور ایستایی سری زمانی، نتایج نرم‌افزار را بهبود می‌بخشد و در نهایت معیارهای آماری قابل‌قبولی برای پیش‌بینی دبی روزانه حاصل شد (Adib et al., 2017). بررسی عملکرد روش برنامه‌ریزی بیان ژن در پیش‌بینی تابش خورشیدی روزانه در گستره ایران بر اساس داده‌های هواشناسی در ۳۱ ایستگاه در گستره ایران نشان داد روش GEP دقیق‌ترین نتایج را در تخمین تابش خورشیدی روزانه در گستره ایران دارد (خادم پور و همکاران، ۱۳۹۷).

سلگی و همکاران (۱۳۹۷) در مدل‌سازی جریان روزانه و ماهانه رودخانه گاماسیاب از مدل برنامه‌ریزی بیان ژن استفاده کردند. برای افزایش عملکرد مدل از دو روش پیش‌پردازش داده‌ها یعنی تبدیل موجک و تجزیه به مؤلفه‌های اصلی PCA استفاده کردند. نتایج نشان داد که عملکرد مدل GEP در دوره ماهانه عملکرد کاهش یافته است. آن‌ها با مقایسه مدل ترکیبی برنامه‌ریزی بیان ژن-موجک با مدل برنامه‌ریزی بیان ژن نشان دادند که عملکرد مدل ترکیبی در هر دو دوره زمانی روزانه و ماهانه از مدل ساده بهتر بوده است. علی‌دادی ده‌کهنه و همکاران (۱۳۹۸) به ارزیابی دو مدل ژنتیکی برنامه‌ریزی ژنتیک و برنامه‌ریزی بیان ژن با استفاده از داده‌های روزانه جریان، دما، بارش و تبخیر در ایستگاه تله‌زنگ جهت مدل‌سازی جریان رودخانه از پرداختن. نتایج نشان داد که مدل برنامه‌ریزی بیان ژن دارای عملکرد بهتری است. علاوه بر این، سرعت اجرای مدل برنامه‌ریزی بیان ژن نسبت به مدل برنامه‌ریزی ژنتیک بیشتر بوده و در زمان کمتری قادر به ارائه نتایج است. همچنین با افزایش تعداد داده‌های ورودی مدل، برنامه‌ریزی ژنتیک کند شده و گاهی قادر به ارائه نتایج نبوده درحالی‌که مدل برنامه‌ریزی بیان ژن این قابلیت را دارد که با تعداد ورودی‌ها و داده‌های بیشتر، نیز عمل مدل‌سازی را انجام دهد، نهایتاً مدل برنامه‌ریزی بیان ژن برای مدل‌سازی و پیش‌بینی جریان رودخانه قابلیت خوبی دارد. عباسی و همکاران (۱۳۹۹) کاربرد برنامه‌ریزی بیان ژن را در پیش‌بینی خشک‌سالی ایستگاه سینوپتیک تبریز موردبررسی قرار دادند و دقت این مدل را قابل‌قبول اعلام کردند. نتایج پژوهش وانگ و همکاران در پیش‌بینی توان تبخیر و تعرق مرجع با GEP دقت زیاد این مدل را نشان می‌دهد. حافظ پرست و رحیمی (۱۳۹۹) به‌منظور ارزیابی تأثیر تغییر

در^۱ GEP، همان نوع از ساختارهای شاخه‌ای مورد استفاده در GP را شامل می‌شود، اما ساختارهای شاخه‌ای که به‌وسیله GEP استنتاج می‌شوند (که بیان درختی نیز نامیده می‌شود) مبین تمامی ژنوم‌های مستقل هستند (Lopes and Weinert, 2004). بنا به اهمیت موضوع، تاکنون پژوهشگران مختلفی در سراسر جهان به مطالعه سری‌های زمانی هیدرولوژیکی از جمله رودخانه‌ها به کمک برنامه‌ریزی بیان ژن پرداخته‌اند (زمانی و همکاران، ۱۳۹۳).

روش برنامه‌ریزی بیان ژن یک روش مناسب و علمی در پیش‌بینی روابط بارش رواناب است. مطالعه پدیده حمل رسوب در آبراهه‌ها و استفاده از روش برنامه‌ریزی ژنتیک یک رهیافت مناسب برای شبیه‌سازی بار معلق است (Aytek and Kisi., 2008). شری و کیس به‌منظور بررسی نوسانات کوتاه‌مدت سطح آب زیرزمینی دو چاه در ترکیه، از روش‌های برنامه‌ریزی بیان ژن و فازی-عصبی استفاده کردند. نتایج به‌دست‌آمده از پژوهش ایشان بیانگر مناسب بودن دو روش بالا در بررسی نوسانات سطح ایستابی بود (Shiri and Kisi., 2011). تراور و گاون از مدل برنامه‌ریزی بیان ژن برای تخمین تبخیر و تعرق در منطقه‌ای در آفریقا پرداختند و دقت این مدل را قابل‌قبول گزارش نمودند (Traore and Guven, 2012).

آزیمتالا و زهیری برای پیش‌بینی دبی جریان در مقاطع مرکب از دو روش برنامه‌ریزی بیان ژن و مدل درختی M5^۲ استفاده کردند. نتایج تحقیقات نشان داد که هرچند هر دو مدل از دقت بالایی برخوردار بودند، اما دقت روش برنامه‌ریزی ژنتیک از مدل درختی M5 بالاتر بود (Azamathulla and Zahiri, 2012).

نتایج بررسی عملکرد روش برنامه‌ریزی بیان ژن را در روند یابی سیلاب رودخانه زنگمار در مقایسه با روش موج دینامیکی نشان داد برنامه‌ریزی بیان ژن قادر است با دقت بیشتری حجم آب نمود خروجی را پیش‌بینی کند. اما الگوی موج دینامیک در خصوص دبی اوج و زمان وقوع آن برتری دارد (قبادیان و همکاران، ۱۳۹۲). نتایج حاصل از بررسی پیش‌بینی عمق سطح آب دریاچه ارومیه حاکی از دقت مطلوب برنامه‌ریزی بیان ژن در شبیه‌سازی نوسانات سطح آب است (کاوه کار و همکاران، ۱۳۹۲). نتایج به‌دست‌آمده از مقایسه مدل‌های فازی عصبی، شبکه عصبی، منحنی سنج و بیان ژن جهت تخمین میزان رسوبات معلق رودخانه آجی چای واقع در آذربایجان شرقی نشان‌دهنده عملکرد بهتر مدل برنامه‌ریزی بیان ژن در مقایسه با سایر مدل‌ها است (ثانی خانی و همکاران، ۱۳۹۴).

امامقلی زاده و همکاران (۱۳۹۵) بارش رواناب حوضه کسلیان را به کمک روش برنامه‌ریزی بیان ژن و با استفاده از پارامترهای هواشناسی و هیدرولوژیکی حوضه، مدل‌سازی کردند. نتایج بیانگر

1- Gene Expression Programming
2- M5 Model Tree

می‌شود. GEP برخلاف GA و GP، چندین عملگر ژنتیکی را برای تکثیر افراد با اصلاحات دارد (Lopes and Weinert., 2004). در GEP، هر ژن به صورت بیان درختی کدگذاری می‌شود. در مورد کروموزوم‌های چندژنی، تمامی بیان درختی‌ها با استفاده از تابع پیوند^۲، از محل گره ریشه^۳ خود به یکدیگر متصل می‌شوند. (Ferreira., 2001).

یک کروموزوم از ژن‌ها تشکیل شده و معمولاً شامل بیش از یک ژن (کروموزوم چند ژن‌ها) است. هر ژن به یکسر^۴ و یک دنباله^۵ تقسیم می‌شود. اندازه رأس (h) به وسیله کاربر تعیین می‌شود، اما اندازه دنباله (t)، به صورت تابعی از h و پارامتر n، به دست می‌آید. پارامتر n، بیشترین تعداد پارامتر مستقل (arity) یافت شده در مجموعه توابع مورد استفاده در اجراست. رابطه (۱)، طول دنباله را با توجه به سایر پارامترها به دست می‌دهد (Lopes and Weinert., 2004).

$$t = h(n-1)+1 \quad (1)$$

به عنوان مثال یک ژن را در نظر بگیرید که از $\{Q, *, /, -, +, a\}$ تشکیل شده است. در این مورد تعداد متغیرهای مستقل $n=2$ است. برای مقادیر طول سر $h=10$ و طول دنباله $t=11$ ، طول ژن برابر با ۲۱ خواهد بود (Ferreira., 2001).

مسائل رگرسیونی نمادین، با استفاده از مجموعه‌ای از توابع و مجموعه‌ای از ترمینال‌ها، مدل‌سازی می‌شوند. مجموعه توابع، معمولاً شامل توابع اصلی حسابی $\{/, *, -, +\}$ ، توابع مثلثاتی یا هر نوع تابع ریاضی دیگر $\{x^2, \exp, \log, \cos, \dots\}$ و یا توابع تعریف شده توسط کاربر است. توابع و ترمینال‌ها، در بخش سر ژن وجود دارند و در قسمت دنباله، فقط ترمینال‌ها وجود دارند (Lopes and Weinert., 2004).

یکی از موارد مهم در GEP، تعیین شاخص اعتبارسنجی تابع هدف است و هدف آن، یافتن راه‌حلی است که برای تمامی موارد برازش به اندازه یک خطای معین به خوبی عمل کند. برخی از توابع برازش مورد استفاده جذر میانگین مربعات خطا (RMSE)، میانگین خطای مطلق (MAE) و ضریب همبستگی (R) است (Lopes et al., 2004). تابع برازش نمایشگر کیفیت یک کروموزوم به عنوان یک جواب از مسئله است و میزان شایستگی یک کروموزوم را در بین جمعیت کروموزومی نشان می‌دهد.

فرایند GEP مانند سایر الگوریتم‌ها زمانی خاتمه می‌یابد که به ضابطه تعیین شده‌ای مثل میزان معین خطای میانگین رسیده باشد و یا چرخه فرایند به تعداد دلخواهی ادامه یافته، در این صورت بهترین راه‌حل یافت شده تا به حال گزارش داده شده و در غیر این صورت

اقلیم بر رواناب منطقه سد جامیشان به بررسی و مقایسه‌ی مدل‌های SVM، GEP و IHACRES پرداختند. نتایج حاصل از پیش‌بینی دبی در هر سه مدل SVM، GEP و IHACRES حاکی از دقت بالاتر مدل‌های IHACRES و GEP نسبت به روش SVM است. یونسو و نوذری (۱۳۹۹) در پیش‌بینی کردن خشک‌سالی کوتاه‌مدت از مدل‌های تلفیقی شبکه عصبی مصنوعی-موجک و برنامه‌ریزی بیان-ژن-موجک استفاده کردند. نتایج پژوهش ایشان نشان داد از میان مدل‌های بررسی شده، برنامه‌ریزی بیان-ژن-موجک با دقت بیش‌تری شاخص بارش معیار و وضعیت خشک‌سالی کوتاه‌مدت را پیش‌بینی می‌کند. همان‌طور که گفته شد، در تحقیقات گذشته پژوهشگران، اثرگذاری پیش‌پردازش (حذف خصوصیات ذاتی داده) به همراه پارامترهای مختلف مدل به صورت هم‌زمان بررسی نشده است. لذا هدف این تحقیق بررسی هم‌زمان خصوصیات داده و پارامترهای مهم در برنامه‌ریزی بیان ژن جهت پیش‌بینی‌های با دقت بیشتر در مراحل آموزش و صحت‌سنجی است.

مواد و روش‌ها

برنامه‌ریزی بیان ژن از جدیدترین الگوریتم‌های فرا کاوشی است که به دلیل دارا بودن دقت کافی، مورد توجه پژوهشگران قرار گرفته است. در این مدل در ابتدای فرایند هیچ‌گونه رابطه تابعی در نظر گرفته نشده و این روش قادر به بهینه‌سازی ساختار مدل و مؤلفه‌های آن است. به طور خلاصه می‌توان گفت در GEP بهسازی‌ها در یک ساختار خطی اتفاق افتاده و سپس به صورت ساختار درختی بیان می‌شود و این موجب می‌شود تنها ژنوم اصلاح شده به نسل بعد منتقل شده و نیازی به ساختارهای سنگین برای تکثیر و جهش وجود نداشته باشد (Ferreira., 2006). به‌طور کلی مراحل اصلی الگوریتم‌های ذکر شده را می‌توان این‌گونه ذکر کرد که ابتدا فرایند با تولید تصادفی کروموزوم‌ها از تعداد معینی افراد (جمعیت اولیه) آغاز می‌شود. سپس در GEP و GP این کروموزوم‌ها به صورت بیان درختی (ET) اظهار می‌شوند. در مرحله بعد با استفاده از تابع هدف میزان برازندگی هر فرد ارزیابی شده و بر اساس میزان عملکردشان انتخاب می‌شوند تا اصلاح شده و فرزندان با ویژگی‌های جدید تشکیل شوند. فرزندان تولید شده، دوباره تحت فرایند توسعه‌ای قرار گرفته تا راه‌حل خوب و مناسبی پیدا شود (Ferreira., 2006). تغییرات ژنتیکی به واسطه عملگرهایی از قبیل جهش، وارون‌سازی، ترانهش IS ترانهش RIS، ترانهش ژن، ترکیب تک نقطه‌ای، ترکیب دونقطه‌ای و ترکیب ژنی است، انجام می‌شود (Ferreira., 2001). روشی که GEP برای انتخاب افراد استفاده می‌کند به روش چرخ رولت^۱ معروف است در این روش اعضا بر اساس میزان سازگاری نسبی آن‌ها انتخاب

2- Linking Function

3- Root Node

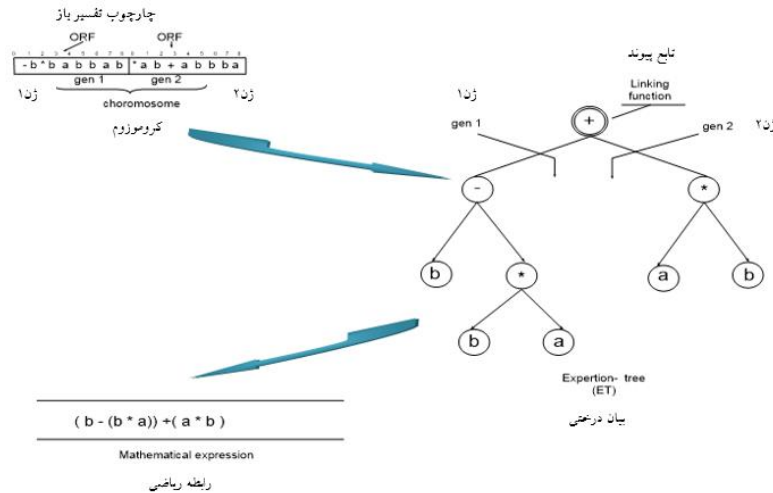
4- Head

5- Tail

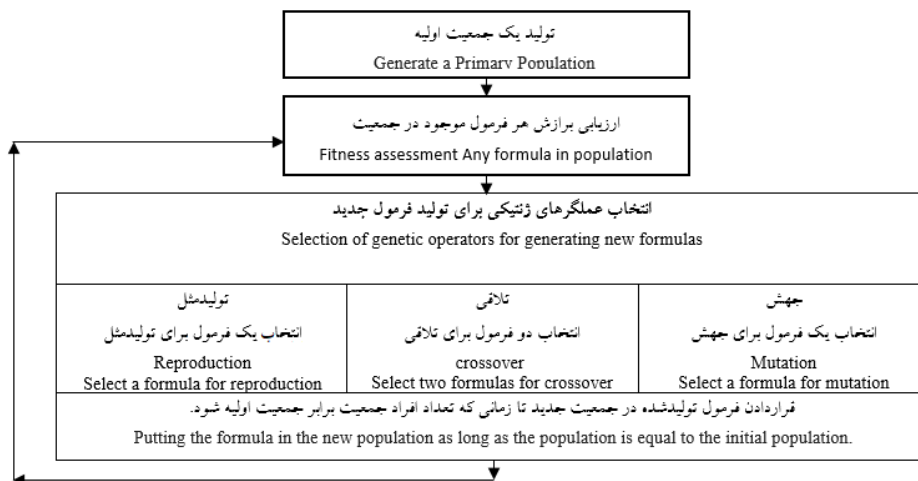
1- Roulette-wheel

بیشتری برای تولید فرزندان دارند.

بهترین راه‌حل از نسل حاضر نگه داشته می‌شود و بقیه راه‌حل‌ها به فرآیندی گزینشی واگذار می‌شوند که بر اساس آن بهترین افراد شانس



شکل ۱- کدگذاری به فرم رشته خطی و کد برداری به صورت ET یک کروموزوم با دو ژن در GEP



شکل ۲- فلوجارت برنامه‌ریزی بیان ژن (Nouryazdan(2015)

مجموعه‌ای شامل ۱۰ برازش ($n=10$) انتخاب گردیده است، بنابراین مقدار تابع برازش برابر با ۱۰۰۰ خواهد بود ($f_{max} = 1000$) مزیت این نوع تابع برازش این است که سیستم می‌تواند با استفاده از آن، راه‌حل بهینه را پیدا کند. جذر میانگین مربعات خطا ($RMSE$)، به‌عنوان معیار خطای تابع برازش انتخاب شده است. کروموزوم‌ها در GEP، معمولاً از چند ژن با طول مساوی تشکیل شده‌اند. برای هر مسئله یا هر اجرا، تعداد ژن‌ها نیز همانند طول رأس قابل انتخاب است. در ادامه در شکل ۲ فلوجارت مدل برنامه‌ریزی بیان ژن و مراحل مختلف الگوریتم آن ارائه شده است.

با تکرار روند ذکر شده و با پیشرفته نسل به جلو کیفیت جمعیت بهبود یافته و به جواب بهینه نزدیک می‌شود (Lopes and Weinert., 2004). از لحاظ ریاضی، برازش f_i از یک برنامه انفرادی i به صورت رابطه (۲) بیان می‌شود.

$$f_i = \sum_{j=1}^n \left(R - \left| \frac{P_{ij} - T_j}{T_j} \right| \cdot 100 \right) \quad (2)$$

که در آن R طول محدوده انتخابی، P_{ij} مقدار پیش‌بینی شده به وسیله برنامه انفرادی i برای مورد برازش j (از میان n مورد برازش) T_j مقدار هدف برای مورد برازش j است. قابل توجه است که عبارت داخل قدر مطلق متناظر با درصد خطای نسبی است (Ferreira., 2006). در این تحقیق طول محدوده انتخابی برابر با ۱۰۰ ($R=100$),

بررسی توابع خودهمبستگی (ACF)

تابع خودهمبستگی یک روش برای بیان وابستگی زمانی در ساختار یک سری زمانی است.

$$\rho_k = \frac{\sum_{i=1}^{n-k} (z_i - \bar{z})(z_{i+k} - \bar{z})}{\sum_{i=1}^n (z_i - \bar{z})^2} \quad -1 \leq \rho_k \leq +1 \quad (3)$$

ρ_k : مقدار تابع خودهمبستگی سری زمانی با تأخیر k

Z_{i+k} و Z_i : مقادیر متغیرها یا داده‌های سری زمانی در مرحله زمانی i و مرحله با تأخیر زمانی k، \bar{z} : مقدار میانگین مربوط به متغیرها (جعفرزاده و همکاران، ۱۳۹۲).

معیارهای ارزیابی

در این تحقیق برای ارزیابی توانایی و دقت مدل برنامه‌ریزی بیان ژن از نمایه‌های ضریب همبستگی (R) و میانگین مربعات خطا (RMSE) و ضریب نش-ساتکلیف (NASH) استفاده شد، که به ترتیب با استفاده از روابط زیر قابل محاسبه است:

$$R = \frac{\sum_{i=1}^n (Q_i - \bar{Q}_i)(\hat{Q}_i - \bar{\hat{Q}}_i)}{\sqrt{\sum_{i=1}^n (Q_i - \bar{Q}_i)^2 \sum_{i=1}^n (\hat{Q}_i - \bar{\hat{Q}}_i)^2}} \quad (4)$$

$$RMS = \sqrt{\frac{\sum_{i=1}^n (Q_i - \hat{Q}_i)^2}{n}} \quad (5)$$

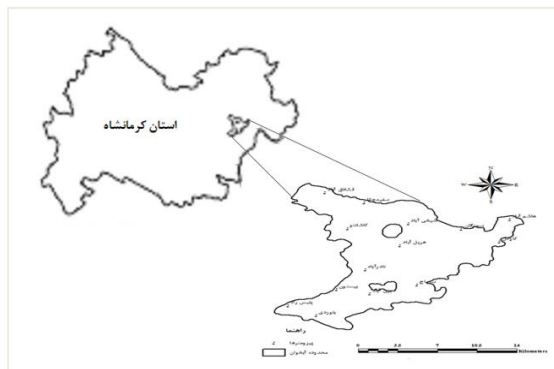
$$NASH = 1 - \frac{\sum_{t=1}^T (\hat{Q}_t^t - Q_t^t)^2}{\sum_{t=1}^T (Q_t^t - \bar{Q}_i)^2} \quad (6)$$

$$NRMSE = \frac{RMSE * 100}{\bar{Q}_i} \quad (7)$$

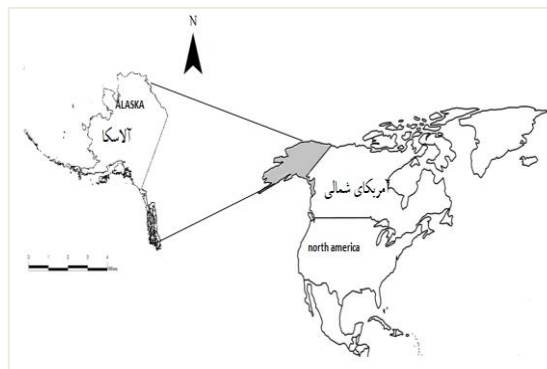
که در آن Q_i مقادیر مشاهده شده، \hat{Q}_i مقادیر محاسبه شده، \bar{Q}_i میانگین مقادیر مشاهده شده و $\bar{\hat{Q}}_i$ میانگین مقادیر محاسبه شده است. علاوه بر معیارهای فوق از نمودارهای مشاهداتی محاسباتی نیز جهت ارزیابی مدل‌ها استفاده گردید. در این تحقیق برای مشاهده وجود روند یا تناوب و نرمال بودن داده‌ها از نرم‌افزار Minitab استفاده شد. همچنین آزمون نرمال بودن داده‌ها از رسم گرافیکی داده‌ها در نمودار توزیع نرمال و با استفاده از آزمون اندرسون-دارلینگ و در سطح اطمینان ۹۵ درصد صورت گرفت. برای به‌کارگیری روش برنامه‌ریزی بیان ژن از نرم‌افزار Genexprotools5.0 و تنظیمات پیش‌فرض برنامه استفاده شد.

منطقه مطالعاتی

برای بررسی اثرگذاری هم‌زمان خصوصیات سری زمانی و پارامترهای تأثیرگذار مدل بر نتایج، از مقادیر اندازه‌گیری شده دمای ماهانه در منطقه آلاسکا واقع در شمال غربی آمریکای شمالی با طول دوره ۵۰ سال و دارای اقلیم سرد و خشک استفاده گردید. تعداد کل داده‌ها در طول دوره (۲۰۱۷-۱۹۶۸) معادل ۶۰۰ داده است. سری زمانی دیگر، داده‌های عمق آب زیرزمینی دشت چمچمال واقع در استان کرمانشاه با طول دوره آماری ۱۲ سال و دارای اقلیم کوهستانی است. از ۷۰ درصد ابتدایی داده‌های مذکور به‌عنوان آموزش و از ۳۰ درصد باقیمانده جهت صحت‌سنجی استفاده گردید.



شکل ۴ - موقعیت جغرافیایی ایستگاه چمچمال



شکل ۳ - موقعیت جغرافیایی آلاسکا

روند ندارند.

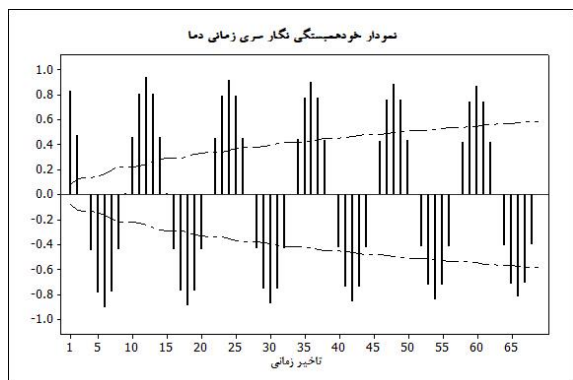
با توجه به شکل (ب) نمودار ACF سری زمانی دمای آلاسکا نشان‌دهنده وجود تناوب سینوسی در این سری زمانی است. این تناوب نامیرا است و طی تأخیرهای مختلف همچنان در سری زمانی وجود دارد. همچنین با توجه به نمودار ACF قدر مطلق بیشترین مقدار تابع

بحث و نتایج

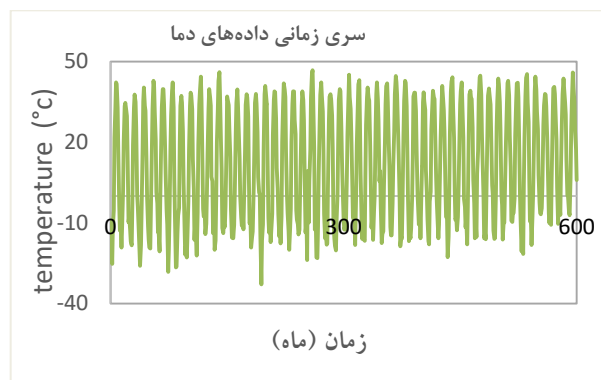
جهت بررسی ایستایی سری زمانی دما، نمودارهای آنالیز حاصل در شکل ۵ (شکل‌های الف تا د) ارائه شده است. همان‌طور که از اشکال ذکر شده مشخص است، داده‌های دما نرمال نیستند، تناوب دارند و با توجه به شیب خط برازش داده‌شده در نمودار روند این داده‌ها

آزمون اندرسون دارلینگ نرمال نیست. با توجه به شکل (د) ناچیز بودن مقدار ضریب t نشان می‌دهد که سری زمانی ماهانه دمای ماهانه فاقد روند است.

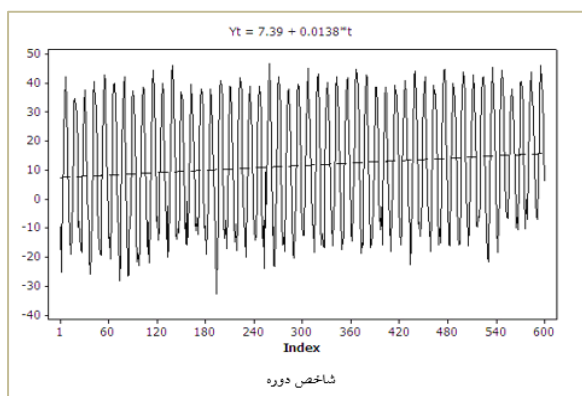
خودهمبستگی سری زمانی دمای آلاسکا نزدیک به یک است و این امر نشان‌دهنده وابستگی زمانی داده‌های سری زمانی است با توجه به شکل (ج)، مقادیر P-Value برای سری زمانی دمای ماهانه کمتر از ۰/۰۵ است. بنابراین این سری زمانی در سطح احتمال ۹۵ درصد و با



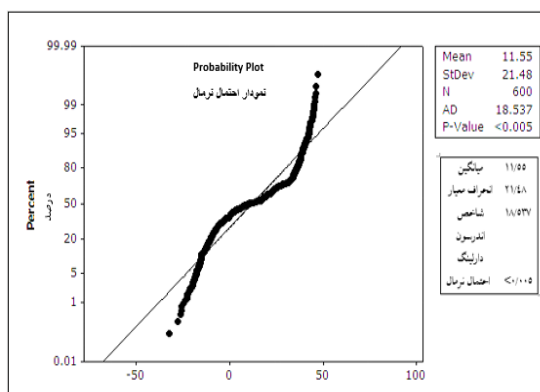
(ب)



(الف)



(د)



(ج)

شکل ۵- نمودارهای سری زمانی داده‌های دما، (الف) نمودار سری زمانی داده‌های دما، (ب) نمودار ACF سری زمانی داده‌های دما، (ج) نمودار نرمال سری زمانی داده‌های دما، (د) نمودار روند سری زمانی داده‌های دما

سه مقدار مختلف در نظر گرفته شد که نتایج اجرای مدل برای هر پارامتر ارائه شده است. در ابتدا از داده‌های ماهانه دمای آلاسکا با دوره زمانی ۵۰ ساله اجرا گرفته شد. تعداد داده‌ها ۶۰۰ ماه است که از ۷۰ درصد داده‌ها برای آموزش مدل و از ۳۰ درصد باقیمانده برای صحت-سنجی مدل استفاده شد. نتایج حاصل در جداول ۲ و ۳ و ۴ ارائه شده است.

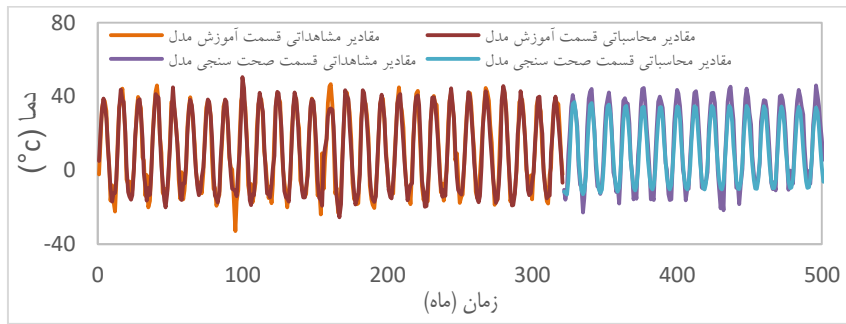
برای اجرای نرم‌افزار بیان ژن پارامترهای مربوط به تعیین عملگرهای ژنتیک و توابع برازش و معیار برازندگی تابع ثابت در نظر گرفته شد و از تنظیمات پیش‌فرض مدل که در جدول ۱ ارائه شده، استفاده شد. برای بررسی تأثیر پارامترهای خود مدل بر نتایج نهایی نرم‌افزار، ساختار کروموزوم‌ها، Embedding Dimension (تعداد تأخیر) و Head Size (اندازه رأس) تغییر داده شده. برای هر پارامتر

جدول ۱- تنظیمات عمومی نرم‌افزار

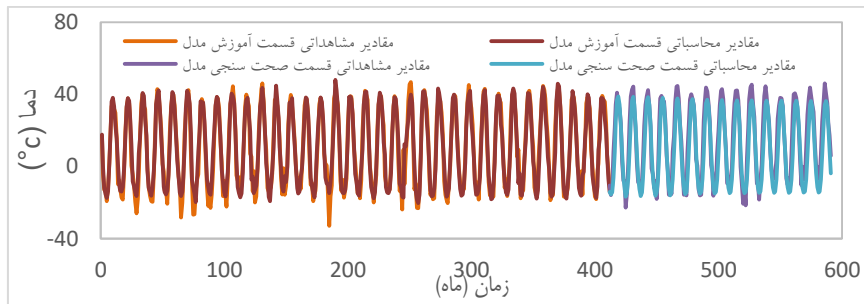
تعداد ژن‌ها در هر کروموزوم	تابع پیوند	نوع تابع برازش	توابع مورد استفاده
۵	Avg2	RMSE	توابع پیش‌فرض سری زمانی

جدول ۲- نتایج بیان ژن برای سری زمانی ۵۰ ساله دمای ماهانه آلاسکا با تعداد تأخیر متغیر در برنامه‌ریزی بیان ژن

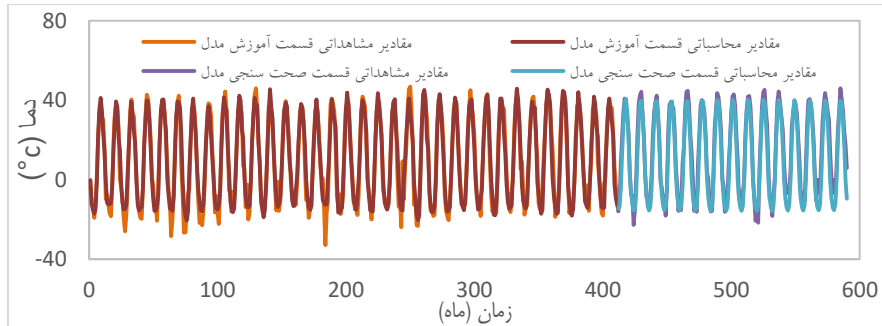
شماره اجرا	تعداد تأخیرها	training			testing				
		R	RMSE	NRMSE	NASH	R	RMSE	NASH	NRMSE
اجرا ۱	۳	۰/۹۷	۸/۲	۷۳/۶	۰/۹۷	۰/۹۵	۶/۷	۰/۹۰	۴۱/۵
اجرا ۲	۱۰	۰/۹۸	۶	۵۲/۰۶	۰/۹۶	۰/۹۷	۵/۵۴	۰/۹۳	۳۸/۳
اجرا ۳	۲۵	۰/۹۷	۵/۸	۵۲/۰۶	۰/۹۵	۰/۹۷	۴/۹	۰/۹۴	۳۳/۹



الف



ب



ج

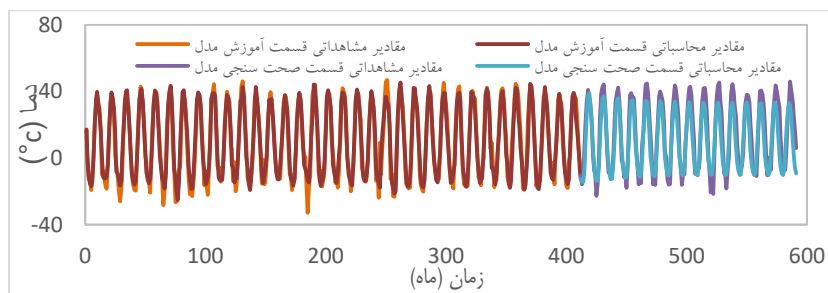
شکل ۶- الف، ب و ج نمودار مشاهداتی محاسباتی دمای ماهانه آلاسکا به ترتیب با تعداد تأخیر ۳، ۱۰ و ۲۵

نداشته است. برای بررسی‌های بیشتر، برای پارامتر اندازه رأس هم همانند مرحله قبل سه مقدار مختلف در نظر گرفته شد و از نرم‌افزار اجرا گرفته شد. نتایج این بخش در جداول ۳ و ۴ و نمودارهای حاصل شده در اشکال ۹ و ۱۰ نشان داده شده است.

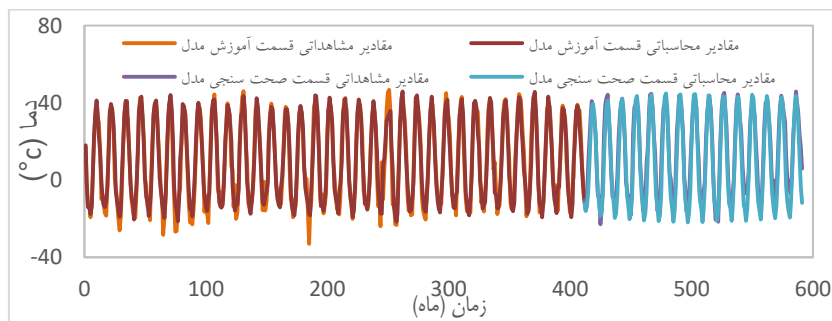
نرم‌افزار در زمان کوتاهی نتایج فوق را حاصل نمود. همان‌طور که از نتایج شکل ۶ پیداست مدل توانسته به خوبی روند تغییرات سری زمانی دما را در هر اجرا پیش‌بینی کند. شاخص R در هر دو قسمت آموزش و صحت‌سنجی مدل برای اجراهای مختلف، در محدوده ۰/۹۷-۰/۹۵ است و تغییر پارامتر تعداد تأخیرها اثر چشم‌گیری بر نتایج

جدول ۳- نتایج بیان ژن برای سری زمانی ۵۰ ساله دمای ماهانه آلاسکا با تعداد اندازه رأس متغیر در برنامه‌ریزی بیان ژن

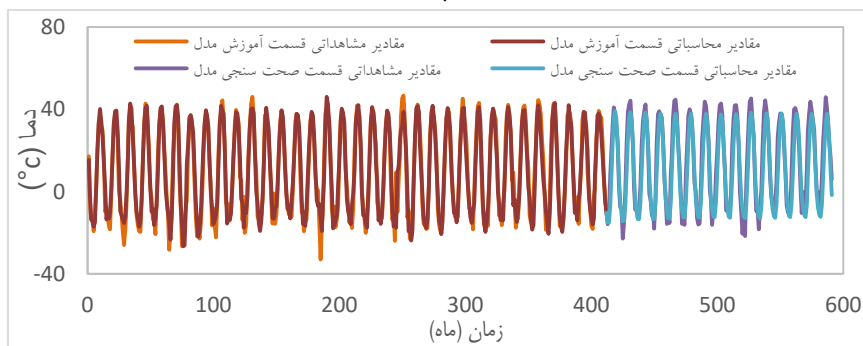
شماره اجرا	اندازه رأس	training				testing			
		R	RMSE	NASH	NRMSE	R	RMSE	NASH	NRMSE
اجرا ۱	۳	۰/۹۵	۱۰	۰/۹۵	۸۹/۷	۰/۹۶	۵/۹	۰/۹۱	۴۰/۸
اجرا ۲	۸	۰/۹۲	۵/۳۸	۰/۹۵	۴۸/۲۹	۰/۹۶	۶/۰۲	۰/۹۴	۴۱/۶
اجرا ۳	۱۵	۰/۹۷	۷/۵	۰/۹۵	۶۷/۳	۰/۹۶	۶/۰۸	۰/۹۲	۴۲/۰



الف



ب



ج

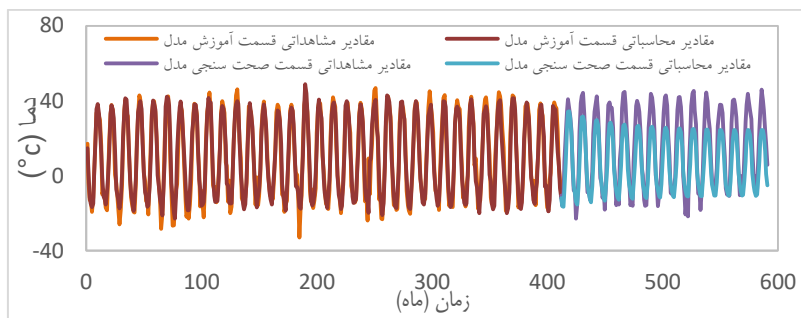
شکل ۹- الف، ب و ج نمودار مشاهداتی محاسباتی دمای ماهانه آلاسکا به ترتیب با تعداد اندازه رأس ۳، ۸ و ۱۵

در هر دو قسمت آموزش و صحت‌سنجی ضرایب به‌دست‌آمده برای معیارهای همبستگی و نش-ساتکلیف بیشتر از ۰/۹۰ هستند.

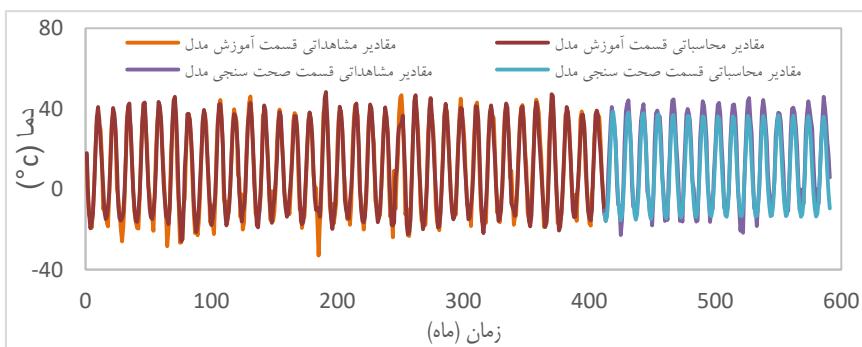
همان‌طور که جدول (۳) نشان می‌دهد، با تعداد رأس متغیر هم نتایج مدل‌سازی توسط برنامه‌ریزی بیان ژن دارای معیارهای آماری قابل قبول در هر دو قسمت آموزش و صحت‌سنجی است. به‌طوری‌که

جدول ۴- نتایج بیان ژن برای سری زمانی ۵۰ ساله دمای ماهانه آلاسکا با تعداد کروموزوم متغیر در برنامه‌ریزی بیان ژن

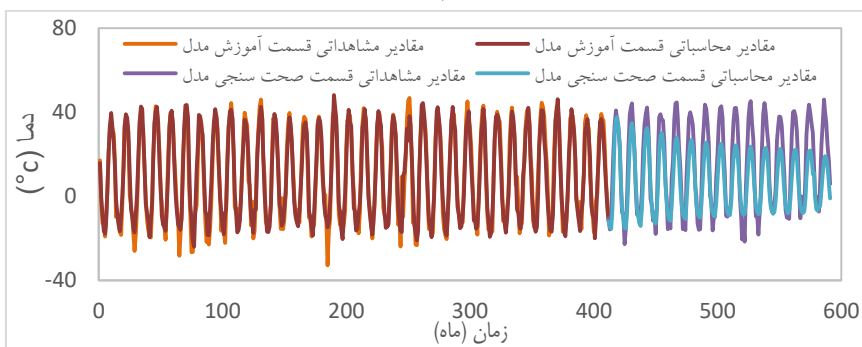
شماره اجرا	تعداد کروموزوم	training				testing			
		R	RMSE	NRMSE	NASH	R	RMSE	NASH	NRMSE
اجرا ۱	۳۰	۰/۹۷	۵/۵	۴۹/۴	۰/۹۹	۰/۹۶	۸/۶	۰/۹۳	۴۷/۷
اجرا ۲	۲۵۰	۰/۹۷	۵/۴	۵۰/۳	۰/۹۹	۰/۹۷	۶/۹	۰/۹۳	۴۷/۷
اجرا ۳	۴۰۰	۰/۹۶	۵/۶	۵۰/۳	۰/۹۹	۰/۹۴	۸/۰۴	۰/۹۱	۴۵/۶



الف



ب

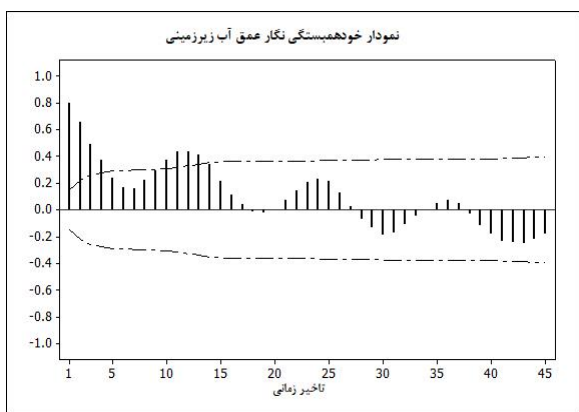


ج

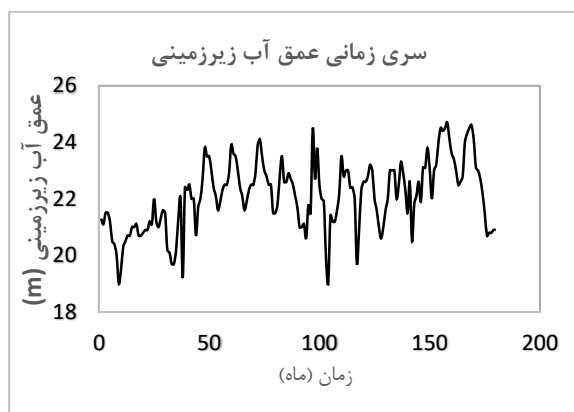
شکل ۱۰- الف، ب و ج نمودار مشاهده‌ای محاسباتی دمای ماهانه آلاسکا به ترتیب با تعداد کروموزوم ۳۰، ۲۵۰ و ۴۰۰

روند ندارند. همچنین با توجه به نمودار خودهمبستگی سری زمانی باگذشت زمان میرا می‌شود. پس می‌توان نتیجه گرفت که سری زمانی ایستا است. انتظار می‌رود که نرم‌افزار برنامه‌ریزی بیان ژن برای داده‌های عمق آب زیرزمینی با توجه به ایستا بودن سری زمانی و با توجه به معیارهای آماری R و $RMSE$ نتایج قابل‌قبول‌تری ارائه دهد. در ادامه این فرضیه بررسی می‌شود. در مرحله دوم داده‌های عمق آب زیرزمینی وارد مدل شد و همانند سری زمانی دما از مدل GEP اجرا گرفته شد. ابتدا تعداد Embedding Dimension (تعداد تأخیر) متغیر در نظر گرفته شد و از سایر تنظیمات پیش‌فرض مدل استفاده شد. در ادامه Head Size و تعداد کروموزوم‌ها متغیر در نظر گرفته شد و از سایر تنظیمات پیش‌فرض مدل استفاده شد. نتایج در جدول‌های شماره ۵ و ۶ و ۷ ارائه شده است.

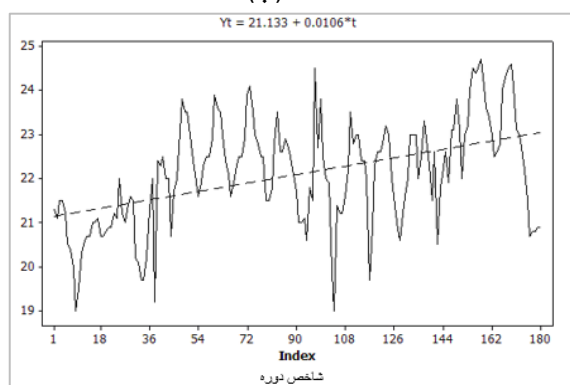
در این مرحله نیز همانند مرحله قبل، نتایج قابل‌قبولی حاصل شده است و نمودارهای مقادیر مشاهده‌شده و محاسبه‌شده در هر دو قسمت آموزش و صحت‌سنجی برهم منطبق هستند. برای بررسی‌های بیشتر، از یک سری زمانی دیگر که مربوط به عمق آب زیرزمینی است، استفاده شد. نمودارهای آنالیز آماری مربوط به این سری زمانی در شکل ۱۱ (نمودارهای الف-د) نشان داده شده است. همان‌طور که شکل (ب) نشان می‌دهد، نمودار سری زمانی داده‌های سطح آب زیرزمینی چمچمال فاقد تناوب سینوسی است و مقادیر حداکثر ناپایداری در نمودار ACF این سری زمانی با تعداد تأخیرهای متفاوت کمتر می‌شود؛ بنابراین نمودار ACF آن نزولی است. داده‌های عمق آب زیرزمینی در سطح ۹۵ درصد اطمینان نرمال هستند، تناوب ندارند و با توجه به شیب خط برازش داده‌شده در نمودار



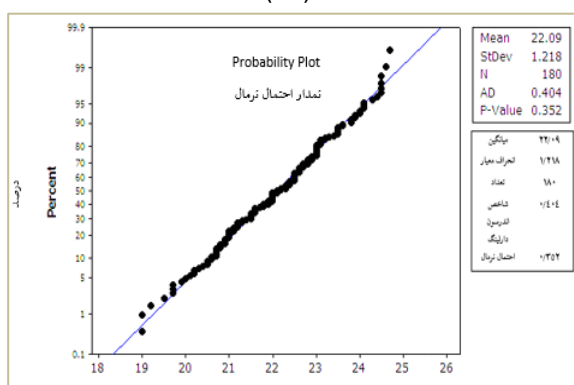
(ب)



(الف)



(د)



(ج)

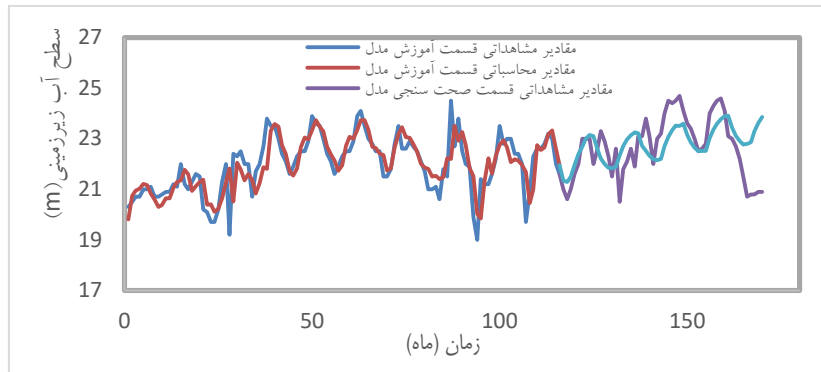
شکل ۱۱- نمودارهای سری زمانی عمق آب زیرزمینی، (الف) نمودار سری زمانی داده‌های آب زیرزمینی، (ب) نمودار ACF سری زمانی داده‌های آب زیرزمینی، (ج) نمودار نرمال سری زمانی داده‌های آب زیرزمینی، (د) نمودار روند سری زمانی داده‌های آب زیرزمینی

جدول ۵- نتایج بیان ژن برای سری زمانی سطح آب زیرزمینی چمچمال با تعداد تأخیر متغیر در برنامه‌ریزی بیان ژن

شماره اجرا	تعداد تأخیرها	training				Testing			
		R	RMSE	NASH	NRMSE	R	RMSE	NASH	NRMSE
اجرا ۱	۳	۰/۷۷	۰/۷۳	۰/۷۰	۳/۳	-۰/۴	۱/۹	-۰/۱	۸/۴
اجرا ۲	۱۰	۰/۷۸	۰/۶۸	۰/۷۲	۳/۱	۰/۴۴	۱/۰۷	۰/۱۵	۴/۷
اجرا ۳	۲۵	۰/۷۳	۰/۷۶	۰/۶۵	۳/۴	۰/۱۴	۱/۵۱	۰/۰۲	۶/۷

همان‌طور که نتایج جدول فوق نشان می‌دهد، مدل برنامه‌ریزی بیان ژن در مدل‌سازی سطح آب زیرزمینی در قسمت آموزش مدل عملکرد بهتری نسبت به قسمت صحت‌سنجی داشته به‌طوری‌که ضریب همبستگی R در محدوده ۰/۷۸-۰/۷۳ حاصل شده است. اما مقدار ضریب همبستگی در قسمت صحت‌سنجی در بیشترین حالت برابر با ۰/۴۴ است؛ بنابراین مدل در پیش‌بینی سطح آب زیرزمینی عملکرد مناسبی ندارد و تغییر در هر پارامتر بیان ژن موجب حاصل شدن ضریب همبستگی بیشتر از ۰/۴۴ در قسمت صحت‌سنجی نشده است. نمودار اجرای شماره ۲ مربوط به قسمت اثرگذاری پارامتر Embedding dimension (تعداد تأخیر) برای هر دو قسمت آموزش و صحت‌سنجی در شکل ۱۲ ارائه شده است. در ادامه تعداد رأس‌ها متغیر در نظر گرفته شد و نرم‌افزار اجرا شد. با توجه به جدول (۶) می‌توان اظهار داشت که تغییر در تعداد رأس هم موجب افزایش ضریب R بیشتر از ۰/۴۲ در قسمت صحت‌سنجی نشده است. در این مرحله هم همانند مرحله قبل عملکرد مدل در قسمت آموزش با حاصل شدن ضریب همبستگی برابر ۰/۸۱ بهتر از قسمت صحت‌سنجی است. در ادامه نمودار اجرای شماره ۳ در شکل ۱۳ ارائه شده است.

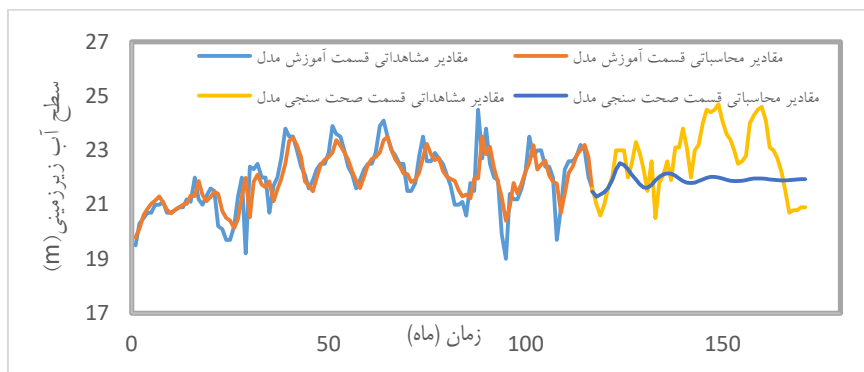
همان‌طور که نتایج جدول فوق نشان می‌دهد، مدل برنامه‌ریزی بیان ژن در مدل‌سازی سطح آب زیرزمینی در قسمت آموزش مدل عملکرد بهتری نسبت به قسمت صحت‌سنجی داشته به‌طوری‌که ضریب همبستگی R در محدوده ۰/۷۸-۰/۷۳ حاصل شده است. اما مقدار ضریب همبستگی در قسمت صحت‌سنجی در بیشترین حالت برابر با ۰/۴۴ است؛ بنابراین مدل در پیش‌بینی سطح آب زیرزمینی عملکرد مناسبی ندارد و تغییر در هر پارامتر بیان ژن موجب حاصل شدن ضریب همبستگی بیشتر از ۰/۴۴ در قسمت صحت‌سنجی نشده است. نمودار اجرای شماره ۲ مربوط به قسمت اثرگذاری پارامتر



شکل ۱۲- نمودار مشاهده‌ای محاسباتی سطح آب زیرزمینی چمچمال با تعداد تأخیر ۱۰

جدول ۶- نتایج بیان ژن برای سری زمانی سطح آب زیرزمینی چمچمال با تعداد اندازه رأس متغیر در برنامه‌ریزی بیان ژن

شماره اجرا	اندازه رأس	training				Testing			
		R	RMSE	NASH	NRMSE	R	RMSE	NASH	NRMSE
اجرا ۱	۳	۰/۷۴	۰/۷۶	۰/۶۵	۳/۵	۰/۳۴	۱/۳۳	۰/۱	۵/۹
اجرا ۲	۸	۰/۶۸	۰/۸۳	۰/۵۵	۳/۸	۰/۱۱	۱/۱۹	۰/۰۱	۵/۲
اجرا ۳	۱۵	۰/۸۱	۰/۶۵	۰/۷۸	۲/۹	۰/۴۲	۱/۲۶	۰/۱۵	۵/۴



شکل ۱۳- نمودار مشاهده‌ای محاسباتی سطح آب زیرزمینی چمچمال با تعداد اندازه رأس ۲۵

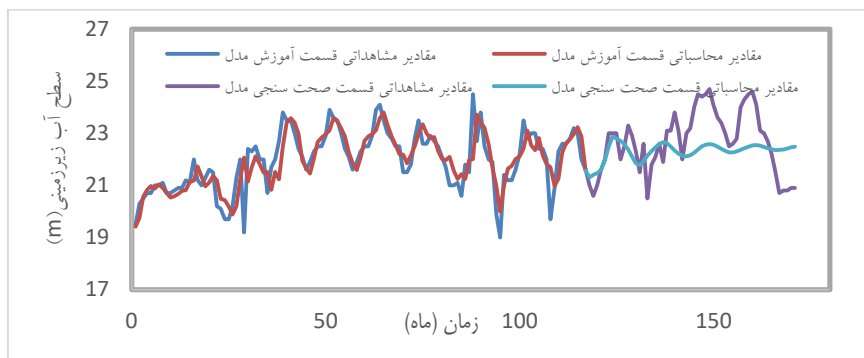
در ادامه تعداد کروموزوم‌ها متغیر در نظر گرفته شد و نرم‌افزار اجرا شد. این نتایج در جدول ۷ ارائه شده است.

جدول ۷- نتایج بیان ژن برای سری زمانی سطح آب زیرزمینی چمچمال با تعداد کروموزوم متغیر در برنامه‌ریزی بیان ژن

شماره اجرا	کروموزوم تعداد	Training				Testing			
		R	RMSE	NASH	NRMSE	R	RMSE	NASH	NRMSE
اجرا ۱	۳۰	۰/۷۸	۰/۶۹	۰/۵۶	۳/۱	۰/۲۸	۱/۲۲	۰/۰۵	۵/۴
اجرا ۲	۲۵۰	۰/۸۲	۰/۶۴	۰/۶۲	۲/۹	۰/۴۴	۱/۳۶	۰/۱۰	۶
اجرا ۳	۴۰۰	۰/۸۰	۰/۶۶	۰/۶۰	۳	۰/۳۳	۱/۳	۰/۱۳	۵/۷

نداشته و ضعیف عملکرده است. اما در قسمت آموزش، مدل برنامه‌ریزی بیان ژن دارای عملکرد بهتر و مناسب‌تری است. دلیل این امر به خاطر پیش‌پردازشی است که روی داده‌ها انجام شده است. در ادامه همانند مراحل قبل با سناریوهای مختلف نرم‌افزار اجرا شد. ابتدا اندازه رأس متغیر در نظر گرفته شد سپس با تعداد کروموزوم‌های متغیر نرم‌افزار اجرا شد. نتایج هر بخش در جداول ۹ و ۱۰ ارائه شده است.

نمودارهای فوق نیز همانند نتایج حاصل شده برای سری زمانی سطح آب زیرزمینی چمچمال در مراحل قبل نیز بیانگر عدم توانایی مدل در تخمین مقادیر داده‌ها در قسمت صحت‌سنجی است. در این حالت با تغییر در هر پارامتر مدل، بهترین نتیجه در قسمت صحت‌سنجی برای معیار R برابر با ۴۴ درصد حاصل شد. همان‌طور که در شکل‌های ۱۲-۱۴ پیداست، نرم‌افزار در مدل‌سازی و پیش‌بینی سری زمانی سطح آب زیرزمینی در قسمت صحت‌سنجی عملکرد قابل‌قبولی



شکل ۱۴- نمودار مشاهداتی محاسباتی سطح آب زیرزمینی چمچمال با تعداد کروموزوم ۲۵۰

داده‌ها باعث افزایش عملکرد مدل GEP شده است. برای پاسخ به فرضیه اثرگذاری خصوصیت داده‌ها بر روی نتایج، به‌منظور حذف تناوب سری زمانی دما عمل استانداردسازی بر روی سری زمانی صورت گرفت. در نتیجه سری زمانی دما نرمال شد و تناوب آن حذف شد و مجدداً این سری زمانی وارد نرم‌افزار شد و همانند مراحل قبل از نرم‌افزار اجرا گرفته شد. نتایج در جداول (۱۰-۸) ارائه شده است.

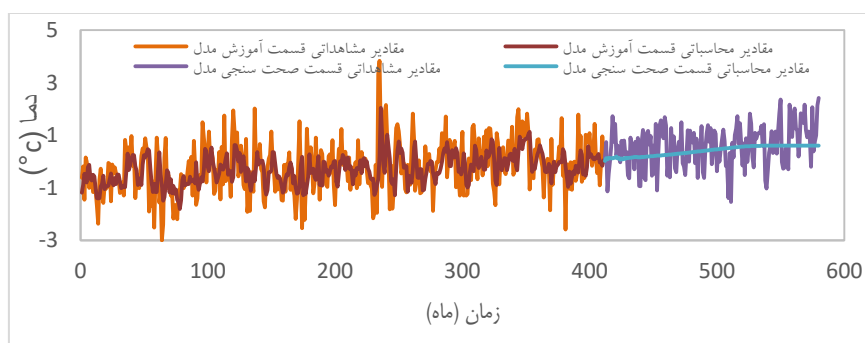
لذا استفاده از روش‌های پیش‌پردازش داده‌ها باعث افزایش عملکرد مدل GEP شده است. سلگی و همکاران (۱۳۹۷) نیز در مدل‌سازی جریان ماهانه رودخانه گاماسیاب از مدل برنامه‌ریزی بیان ژن استفاده کردند. نتایج آن‌ها نشان داد که عملکرد مدل ترکیبی از مدل ساده بهتر بوده است. دلیل این امر به خاطر پیش‌پردازی است که روی داده‌ها پیاده شده بود. لذا استفاده از روش‌های پیش‌پردازش

جدول ۸- نتایج بیان ژن برای سری زمانی استاندارد شده ۵۰ ساله دمای آلاسکا با تعداد تأخیر متغیر در برنامه‌ریزی بیان ژن

شماره اجرا	تعداد تأخیرها	training			testing		
		R	RMSE	NASH	R	RMSE	NASH
اجرا ۱	۳	-۰/۴۲	-۰/۸۷	-۰/۱۳	-۰/۱	۱/۱۳	-۱/۰۲
اجرا ۲	۱۰	-۰/۵	-۰/۸۵	-۰/۱۸	-۰/۱	۱/۳۲	-۱/۸
اجرا ۳	۲۵	-۰/۳۴	-۰/۹۴	۰	۰/۲۶	۰/۷۸	۰/۰۲

۰/۵ یا ۵۰٪ و برای قسمت صحت‌سنجی برابر ۰/۲۶ حاصل شده است. نمودار اجرای شماره ۳ در این بخش در شکل نشان داده شده است.

همان‌طور که جدول فوق نشان می‌دهد، بعد از اعمال استانداردسازی سری زمانی دما نتایج مدل‌سازی کاهش چشم‌گیری داشته و معیارهای آماری قابل قبولی حاصل نشده است. به‌طوری‌که بهترین نتیجه ضریب همبستگی (R) برای قسمت آموزش مدل برابر



شکل ۱۵- نمودار مشاهداتی محاسباتی دمای ماهانه استاندارد شده آلاسکا با تعداد تأخیر ۲۵

جدول ۹- نتایج بیان ژن برای سری زمانی استاندارد شده ۵۰ ساله دمای آلاسکا با تعداد اندازه رأس متغیر در برنامه ریزی بیان ژن

شماره اجرا	اندازه رأس	training			testing		
		R	RMSE	NASH	R	RMSE	NASH
اجرا ۱	۳	-۰/۱	۱/۵۲	-۲/۶	۰/۴۸	۰/۸۵	۰/۱۵
اجرا ۲	۸	۰/۱۱	۱/۲۸	-۱/۶۱	۰/۵۲	۰/۸۴	۰/۲۱
اجرا ۳	۱۵	۰/۱۱	۱/۲۷	-۱/۵	۰/۵	۰/۸۴	۰/۲

مدل ۰/۱۱ است و در قسمت صحت‌سنجی برای معیار R مقدار بیشتر از ۰/۵۲ حاصل نشده است.

همان‌طور که جدول (۹) نشان می‌دهد، نتایج مدل‌سازی سری زمانی استاندارد شده دمای آلاسکا به طرز معنی‌داری از لحاظ معیارهای کاهش یافته است. به طوری که بیشترین مقدار R در قسمت آموزش

جدول ۱۰- نتایج بیان ژن برای سری زمانی استاندارد شده ۵۰ ساله دمای آلاسکا با تعداد کروموزوم متغیر در برنامه ریزی بیان ژن

شماره اجرا	تعداد کروموزوم‌ها	training			testing		
		R	RMSE	NASH	R	RMSE	NASH
اجرا ۱	۳۰	۰/۱۹	۰/۹۸	-۰/۰۷	-۰/۲۱	۰/۹۸	-۰/۵
اجرا ۲	۲۵۰	۰/۱۲	۰/۹۷	-۰/۰۲	-۰/۰۷	۱/۵	-۱/۱
اجرا ۳	۴۰۰	۰/۵	۰/۸۵	۰/۲	-۰/۱۲	۱/۲۷	-۱/۵

پارامترهای تأثیرگذار خود نرم‌افزار شامل Embedding Dimension (تعداد تأخیر)، Head Size (اندازه رأس)، و اندازه طول کروموزوم و تأثیر هم‌زمان آن‌ها در پیش‌بینی نتایج بهتر مورد بررسی قرار گرفت. به همین منظور ابتدا از سری زمانی دمای ماهانه منطقه آلاسکا که دارای خصوصیت قطعی تناوب است، استفاده شد. با تغییر هر پارامتر تأثیرگذار بیان ژن بهترین نتایج حاصل شده و تغییر در پارامترهای مدل تأثیر معنی‌داری در نتایج نداشت. نتایج به دست آمده برای سری زمانی دما نشان دهنده عملکرد مناسب و قابل قبول و مورد تأیید این الگوریتم از لحاظ معیارهای آماری است. همان‌طور که در شکل‌های شماره (۱۰-۶) مشخص است، نمودارهای برازش داده شده توسط مدل بر نمودار خود داده‌ها منطبق است و نمودارهای مشاهداتی و محاسباتی دارای هم‌پوشانی خوبی باهم هستند. به طوری که مقدار ضریب همبستگی (R) در هر اجرا با هر پارامتر تأثیرگذار بیان ژن برای هر دو قسمت آموزش و صحت‌سنجی در محدوده ۰/۹۷-۰/۹۰ است. علت این امر وجود ترم تناوب که از خصوصیات داده‌های ماهانه است، است. برای اثبات این موضوع با استفاده از استانداردسازی این ویژگی از داده‌های ماهانه حذف شد و با تغییر در هر پارامتر بیان ژن نتایج معیار R بهتر از ۵۲ درصد نشد. نمودارهای مشاهداتی و محاسباتی در شکل ۱۵ نشان می‌دهد که مدل در این حالت قادر نبوده مقادیر را تخمین بزند. در مرحله بعد از داده‌های ایستگاه چمچمال که فاقد خصوصیت تناوبی بود و ایستا بود، استفاده شد. در این حالت با تغییر در هر پارامتر مدل، بهترین نتیجه در قسمت صحت‌سنجی برای معیار R برابر با ۴۴ درصد حاصل شد. همان‌طور که در شکل‌های

همان‌طور که جدول (۱۰) نشان می‌دهد، نتایج مدل‌سازی سری زمانی استاندارد شده دمای آلاسکا به طرز معنی‌داری از لحاظ معیارهای آماری کاهش یافته است. به طوری که بیشترین مقدار R در قسمت آموزش مدل ۰/۱۹ است و در قسمت صحت‌سنجی برای معیار R مقدار منفی به دست آمده است. نتایج جدول فوق نشان می‌دهد که نرم‌افزار قادر به مدل‌سازی سری زمانی استاندارد شده دمای آلاسکا نبوده و معیارهای آماری قابل قبول حاصل نشده است. اعمال پیش-پردازش بر روی سری زمانی دما و حذف خصوصیت تناوبی آن، باعث کاهش معنی‌دار معیارهای آماری گردید و عملکرد مدل به طرز چشم-گیری کاهش یافت. موثنی و همکاران بر غالب بودن ترم تناوب و اثر آن در بالا بودن عملکرد نتایج پیش‌بینی مدل‌های ARIMA بر سری زمانی دما اشاره داشتند. آن‌ها در مطالعات خود از شش سری زمانی شامل ۴ سری جریان و ۲ سری دمای آب استفاده نمودند. برای تجزیه و تحلیل داده‌ها ۹۲۲۸ مدل معرفی اجرا گردیده و نتایج آن‌ها نشان داده است که استانداردسازی بهترین روش برای حذف تناوب سری زمانی دمای ماهانه آب و مهم‌ترین عامل در دقت مدل‌های پیش‌بینی است (Moeeni et al., 2017)

نتیجه‌گیری

در این تحقیق برای پاسخ به فرضیه تأثیر خصوصیات مربوط به یک سری زمانی بر روی نتایج مدل‌سازی، از الگوریتم برنامه‌ریزی بیان ژن و نرم‌افزار GeneXprotools استفاده گردید و از دونقطه نظر تأثیر پیش‌پردازش سری زمانی شامل تناوب، روند، و نرمال بودن و

بردار پشتیبان، نشریه تحقیقات کاربردی علوم جغرافیایی. ۱۸ (۵۰): ۰۹۱-۱۰۳.

ظهیری، ع. و قربانی، خ. ۱۳۹۲. شبیه‌سازی دبی جریان در مقاطع مرکب به کمک مدل درخت تصمیم M5. نشریه پژوهش‌های حفاظت آب‌و خاک. ۲۰ (۳): ۱۱۳-۱۳۲.

عباسی، ع.، خلیلی، ک.، بهمنش، ج. و شیرزاد، ا. ۱۳۹۹. کاربرد برنامه‌ریزی بیان ژن در پیش‌بینی خشک‌سالی (مطالعه موردی: ایستگاه سینوپتیک تبریز). مجله محیط‌زیست و مهندسی آب. ۵ (۱): ۱-۱۴.

علی دادی ده کهنه، ص.، سلگی، ا.، شهینی دارابی، م. و زارعی، ح. ۱۳۹۸. ارزیابی مدل‌های ژنتیکی جهت مدل‌سازی جریان رودخانه. مهندسی آبیاری و آب ایران. ۹ (۳۵): ۱-۱۷.

قبادیان، ر.، قربانی، ع.م. و خلج، م. ۱۳۹۲. بررسی عملکرد روش برنامه‌ریزی بیان ژن در روند یابی سیلاب رودخانه زنگمار در مقایسه با روش موج دینامیکی. نشریه آب‌و خاک (علوم و صنایع کشاورزی) ۲۷ (۳): ۵۹۲-۶۰۲.

کاوکار، ش.، نعمتی، س. و ازانی، ع. ۱۳۹۲. ساختار درختی برنامه‌ریزی بیان ژن جهت شبیه‌سازی نوسانات تراز آب، کنفرانس بین‌المللی عمران، معماری و توسعه پایدار شهری، تبریز، دانشگاه آزاد اسلامی واحد تبریز.

مهدی زاده، س.، ۱۳۹۶. مدل‌سازی و پیش‌بینی برخی از پارامترهای مورد استفاده در مدیریت آب و کشاورزی با استفاده از مدل‌های هوش مصنوعی و تلفیق مدل‌های مذکور با روش‌های استوکاستیک. رساله دکتری تخصصی. دانشگاه ارومیه.

یونسی، م. و نوذری، ح. ۱۳۹۹. ارزیابی مدل‌های تلفیقی شبکه‌ی عصبی مصنوعی-موجک و برنامه‌ریزی بیان ژن-موجک در پیش‌بینی کردن خشک‌سالی کوتاه‌مدت. نشریه پژوهش‌های آبخیزداری (پژوهش و سازندگی). ۳۳ (۱): ۳۹-۵۵.

Adib A., Mahmoudian Kafshgar Kalae M., Mahmoudian Shoushtari M., Khalili K. 2017. Using of gene expression programming and climatic data for forecasting flow discharge by considering trend, normality, and stationarity analysis. *Arabian Journal of Geosciences* 10(9): 208.

Aytek A., Kisi O. 2008. A genetic programming approach to suspended sediment modeling. *Journal of Hydrology* 351(3-4): 288-298.

Azamathulla H.Md., Zahiri A. 2012. Flow discharge prediction in compound channels using linear genetic programming. *Journal of Hydrology*. 454-455(6): 203-207.

(۱۴-۱۲) پیداست، نرم‌افزار در مدل‌سازی سری عمق آب زیرزمینی، عملکرد قابل‌قبولی نداشته و ضعیف عملکردده است.

در نهایت می‌توان نتیجه گرفت تأثیر پیش‌پردازش سری زمانی در نتایج پیش‌بینی با مدل برنامه‌ریزی بیان ژن خیلی بیشتر از پارامترهای خود مدل است. دارا بودن مؤلفه تناوب برای سری زمانی موجب عملکرد عالی مدل می‌شود. همچنین در مورد نتایج به‌دست‌آمده برای سری زمانی عمق آب زیرزمینی می‌توان اظهار داشت که تغییر در پارامترهای خود نرم‌افزار تنها درصد کمی حدود ۱۰ درصد در بهبود نتایج تأثیر دارد. و تنها خصوصیات خود داده‌ها است که در عملکرد مدل اثرگذار است.

منابع

امامقلی زاده، ص.، کریمی‌دمنه، ر. و مهدی پناه، ح. ۱۳۹۵. برآورد رواناب حوضه آبریز کسلیان با استفاده از روش برنامه‌ریزی بیان ژن. *مجله مطالعات منابع طبیعی، محیط‌زیست و کشاورزی*. ۲ (۴): ۱-۷.

ثانی خانی، ه.، نیک پور، م.، فرسادی‌زاده، د. و معیری م.م. ۱۳۹۴. پیش‌بینی بار معلق رودخانه با استفاده از سامانه‌های هوشمند. *مجله پژوهش آب ایران*. ۹ (۲): ۱۶۵-۱۶۸.

جعفر زاده، ج.، رستم زاده، ه. و اسدی، ا. ۱۳۹۶. مدل‌سازی زمانی تراز آب زیرزمینی با استفاده از روش‌های پایه تحلیل سری‌های زمانی (مطالعه موردی: دشت اردبیل). *نشریه دانش آب‌و خاک*. ۲۷ (۴): ۱۸۵-۱۹۶.

حافظ پرست مودت، م. و رحیمی، ب. ۱۳۹۹. مقایسه مدل‌های SVM، GEP و IHACRES در پیش‌بینی تغییرات رواناب ناشی از تغییر اقلیم (مطالعه موردی: سد جامیشان). *نشریه تحقیقات آب‌و خاک ایران*. ۵۹ (۱۰): ۲۴۸۳-۲۴۹۹.

خادم پور، ف.، خاشعی سیوکی، ع. و امیرآبادی زاده، م. ۱۳۹۷. بررسی عملکرد روش برنامه‌ریزی بیان ژن در پیش‌بینی تابش خورشیدی روزانه در گستره ایران. *نشریه پژوهش‌های اقلیم‌شناسی*. ۹ (۳۶): ۴۳-۵۶.

زمانی، ر.، احمدی، ف. و رادمنش، ف.، ۱۳۹۳. مقایسه روش‌های برنامه‌ریزی بیان ژن، سری زمانی غیرخطی، خطی و شبکه عصبی مصنوعی در تخمین دبی روزانه (مطالعه موردی: رودخانه کارون)، *نشریه آب‌و خاک (علوم و صنایع کشاورزی)*. ۲۸ (۶): ۱۱۷۲-۱۱۸۲.

سلگی، ا.، زارعی، ح.، دارابی، م. و ده کهنه، ص. ۱۳۹۷. پیش‌بینی بارش ماهانه با استفاده از مدل‌های برنامه‌ریزی بیان ژن و ماشین

- Journal of Hydrology. 547: 348-364.
- Salas J. D., Delleur j.w, Yevjevich v. and Lane,w.l. 1996. Applied Time Series in Hydrology, Mc Graw Hill.
- Shiri J., Kisi O. 2011. Comparison of genetic programming with neuro-fuzzy systems for predicting short-term water table depth fluctuations. Comput. Geosci. 37: 1692-1701.
- Traore S., Guven A. 2012. Regional-specific numerical models of evapotranspiration using gene-expression programming interface in Sahel. Water Resources. Management. 26: 4367-4380.
- Wang Sh., Lian J., Peng Y., Hu B., Chen H. 2019. Generalized reference evapotranspiration models with limited climatic data based on random forest and gene expression programming in Guangxi, China. Agricultural Water Management. 221: 220-230.
- Ferreira C. 2001. Gene expression programming: a new adaptive algorithm for solving problems.complex Systems.13 (2):87-129.
- Ferreira C. 2006. Gene expression programming: Mathematical Modeling by an Artificial Intelligence. 2nd edn. Springer-Verlag Germany.463 p.
- Lopes H.S., Weinert W.R. 2004. EGIPSYS: An enhanced gene expression programming approach for symbolic regression problems. International Journal of Applied Mathematics and Computer Science, 14(3): 375-384.
- Moeeni H., Bonakdari H., Fatemi S.E., Zaji A.H. 2017. Assessment of Stochastic Models and a Hybrid Artificial Neural Network-Genetic Algorithm Method in Forecasting Monthly Reservoir Inflow. INAE Letters. 2(1): 13-23.
- Moeeni, H., Bonakdari H., Fatemi, S.E, 2017. Stochastic model stationarization by eliminating the periodic term and its effect on time series prediction.

Assessment Effects of Data Preprocessing and Modeling Parameters of Gene Expression Programming on Accuracy of Time Series Forecasting

M. Salehi^{1*}, S.E. Fatemi²

Received: Jan.19, 2021

Accepted: Apr.13, 2021

Abstract

Hydrological time-series is a time-dependent hydrological variable that finding the model of changes and predicting is the most important goal of time-series analysis. The purpose of this study is to simultaneously study the characteristics of time series and their prediction and the important parameters of the GEP for high-precision predictions in the training and validation. In this study, groundwater depth time-series of Chamchamal plain station located in Kermanshah province with a 12-year period and mountainous climate and the monthly time-series of Alaska temperature with a 50-year period and cold and dry climate have been used. Genexprotools5.0 software has been used to model time-series by GEP.

The results of studying with GEP showed that the periodicity of data properties that existed in the time series of temperature caused correlation results above 90% in different stages of training and validation. So that the effect of different parameters of GEP is less than 10% in improving results. On the other hand, by examining the time-series of groundwater depth, which lacks periodicity and has a descending ACF shape, the prediction results of the GEP with any effective expression parameter, R more than 44% in the validation wasn't obtained. This means that the time-series preprocessing has a greater impact on the prediction results. So that by eliminating the semester, the prediction results in all stages of modeling are significantly reduced. In this case, the best R for the validation is 50%.

Keywords: Forecasting, Gene Expression, Periodicity, Preprocessing, Time Series

1- Master of Water Resources, Water Engineering Department, Campus of Agricultur and Natural Resources, Razi University, Iran

2- Assistant Professor, Water Engineering Department, Campus of Agricultur and Natural Resources, Razi University, Iran

(*- Corresponding Author Email: e_fatemi78@yahoo.com)