

مقایسه‌ی مدل‌های درخت تصمیم M5 و الگوریتم نزدیک‌ترین همسایگی K در پیش‌بینی بارش ماهانه (مطالعه موردی: ایستگاه سینوپتیک بیرجند)

فاطمه پورصالحی^{۱*}، علی شهیدی^۲، عباس خاشعی سیوکی^۳

تاریخ دریافت: ۱۳۹۷/۱۲/۱۷ تاریخ پذیرش: ۱۳۹۸/۴/۳

چکیده

باتوجه به واقع شدن ایران در اقلیم خشک و نیمه خشک، توزیع ناهمگن بارندگی و همچنین وقوع پدیده‌ی تغییر اقلیم سبب ایجاد پدیده‌هایی مانند سیل، خشکسالی، بیابان‌زایی و تولید ریزگردها و نیز ایجاد خسارت‌های مختلف اقتصادی، اجتماعی و زیست‌محیطی گردیده است. یکی از راهکارهای اولیه جهت کاهش این خسارات، پیش‌بینی رخداد بارندگی است. هدف از مطالعه‌ی حاضر پیش‌بینی بارش ماهانه با بکارگیری روش‌های داده‌کاوی الگوریتم‌های درخت تصمیم (M5) و نزدیک‌ترین همسایگی K (KNN) و مقایسه‌ی این دو روش در راستای تعیین روش کارتر در زمینه‌ی پیش‌بینی بارندگی با استفاده از داده‌های هواشناسی ماهانه‌ی ایستگاه سینوپتیک بیرجند طی دوره‌ی آماری ۲۰۱۰-۱۹۶۱ میلادی در سه حالت داده خام، میانگین متحرک سه‌ساله و میانگین متحرک پنج‌ساله در نرم افزار Weka می‌باشد. نتایج نشان داد که در تمامی سناریوهای تعریف شده، مدل درختی M5 نسبت به مدل KNN توانایی بیشتری در پیش‌بینی بارش ماهانه‌ی این ایستگاه دارد. همچنین پس از بررسی معیارهای ارزیابی R، RMSE، MAE و NS، سناریو پانزدهم با پارامترهای ورودی اختلاف میانگین حداکثر و حداقل دما، متوسط رطوبت نسبی، میانگین سرعت باد و درجه روز سرمایش (بر پایه ۲۱ درجه سانتی‌گراد) در هر ماه به عنوان بهترین سناریو برای پیش‌بینی بارش همان ماه تعیین گردید. همچنین نتایج به دست آمده از مقایسه‌ی سناریوهای تعریف شده در هر مدل در سه حالت داده‌های خام، میانگین متحرک سه‌ساله و میانگین متحرک پنج‌ساله نشان می‌دهد که در اکثر سناریوها میانگین متحرک پنج‌ساله به طور میانگین با مقادیر $R=0/904$ ، $RMSE=6/054$ و $MAE=4/780$ در مدل M5 و به طور میانگین با مقادیر $R=0/837$ ، $RMSE=7/698$ و $MAE=5/595$ در مدل KNN پیش‌بینی دقیق‌تری از بارش ماهانه را ارائه می‌دهد.

واژه‌های کلیدی: خشکسالی، درخت تصمیم، روش‌های داده‌کاوی، نرم‌افزار Weka، نزدیک‌ترین همسایگی K

مقدمه

ویژه‌ای جهت بهینه‌سازی صرف هزینه‌ها و استفاده از این منابع برخوردار است. بارش، پدیده‌ای بسیار پیچیده، غیرخطی و نسبت به زمان و مکان متغیر می‌باشد. عوامل مؤثر زیادی در تغییرات آن نقش دارند؛ که به طور کلی می‌توان آن‌ها را به دو دسته‌ی اقلیمی و جغرافیایی تقسیم کرد. از جمله عوامل اقلیمی می‌توان به رطوبت، فشار، دما، پوشش ابر، سرعت باد، عوامل فصلی و همچنین سیگنال‌های بزرگ مقیاس اقلیمی و به خصوص دوره‌های النینو و لانینا اشاره نمود. از عوامل جغرافیایی نیز می‌توان به دوری و نزدیکی از مراکز تولید جبهه (مثلاً دریاها و یا صحراهای بزرگ) و همچنین ارتفاع، اشاره کرد (مهدوی، ۱۳۷۴).

بنابراین با توجه به اهمیت بارش به عنوان یکی از پارامترهای مهم هواشناسی، شناخت کافی از مقدار و همچنین عوامل مؤثر بر این عنصر، بررسی تغییرات و پیش‌بینی آن، برای رسیدن به برنامه‌ریزی کارتر در مدیریت بخش‌های کشاورزی، اقتصادی و اجتماعی و همچنین پیش‌بینی پدیده‌های هیدرولوژیکی مانند سیلاب و خشکسالی ضروری می‌باشد. در علم آمار روش‌های مختلفی برای

کشور ایران بر روی کمربند خشک جهان قرار دارد و با بارندگی معادل یک سوم متوسط جهانی، کشوری خشک است و به این دلیل، خشکی جزء صفات ذاتی آن محسوب می‌شود. روند بارندگی در ایران حاکی از آن است که این کشور به سوی خشکی پیش می‌رود و می‌بایست برنامه‌ریزی‌ها و تدابیر در مدیریت منابع آب بر این اساس پی‌ریزی شود. نتایج مطالعه و بررسی بلند مدت بارندگی و درجه حرارت نشان می‌دهد که در مجموع، ۶۵ درصد از اراضی کشور در گستره‌ی اقلیم‌های خشک و فراخشک قرار دارند (نعیمی و احقافی، ۱۳۸۱). پیش‌بینی بارش و برآورد نزولات جوی، به عنوان یکی از مهم‌ترین پارامترهای اقلیمی در حوزه‌ی مدیریت منابع آب، از اهمیت

۱- دانشجوی دکتری مهندسی منابع آب دانشگاه بیرجند

۲- دانشیار گروه علوم و مهندسی آب دانشگاه بیرجند

۳- دانشیار گروه علوم و مهندسی آب دانشگاه بیرجند

* - نویسنده مسئول: (Email: Fatemehpursalehi@birjand.ac.ir)

مشاهداتی هستند و مقادیر حدی را نمی‌توانند به خوبی پیش‌بینی کنند. با این وجود روش نزدیکترین همسایگی نسبت به دیگر روش‌ها نتایج بهتری را ارائه کرد. امیدوار و همکاران (۱۳۹۳)، توانایی مدل درخت تصمیم را در پیش‌بینی بارش ایستگاه سینوپتیک کرمانشاه مورد ارزیابی قرار دادند. نتایج نشان داد که در ایستگاه سینوپتیک کرمانشاه درخت تصمیم‌گیری رگرسیونی، مدلی نسبتاً کارا در پیش‌بینی بارش می‌باشد که استفاده از میانگین متحرک نسبت به سایر حالات منجر به افزایش چشمگیر کارایی مدل درخت تصمیم می‌شود و در صورت تعدیل دامنه تغییرات داده‌های ورودی قادر است با ضریب اطمینان بالایی میزان بارش را ۳۰ ماه قبل از وقوع برآورد نماید که در شبیه‌سازی‌های صورت گرفته، زمانی که از میانگین متحرک پنج ساله داده‌ها برای اجرای مدل استفاده گردیده، ترکیب بارش قبلی، دمای ماکزیمم به عنوان مناسب‌ترین حالت شناسایی شده است. دستورانی و همکاران (۱۳۹۱)، کارایی مدل درخت تصمیم را در پیش‌بینی بارش ایستگاه سینوپتیک یزد تعیین نمودند. نتایج نشان داد که در ایستگاه یزد، مدل درخت تصمیم‌گیری خصوصاً در شرایطی که از میانگین متحرک ۵ ساله داده‌ها استفاده گردد، دارای توانایی مناسبی در پیش‌بینی میزان بارش می‌باشد. ستاری و نهرین (۱۳۹۲)، با استفاده از سیستم‌های هوشمند (شبکه‌ی عصبی مصنوعی و برنامه‌ریزی ژنتیک) مقادیر حداکثر بارش روزانه‌ی ایستگاه‌های اهر و جلفا را پیش‌بینی و با مدل درختی M5 مقایسه نمودند. در حالت کلی می‌توان گفت که هر سه روش مذکور ضمن رقابت با یکدیگر نتایج نسبتاً دقیقی را جهت پیش‌بینی حداکثر بارش روزانه در ماه مورد نظر در منطقه ارائه می‌کنند ولی به دلیل ارائه روابط خطی ساده و قابل فهم توسط مدل درختی M5، این روش می‌تواند به عنوان روشی کاربردی و جایگزین برای محاسبه حداکثر بارش روزانه در ماه مورد توجه قرار گیرد.

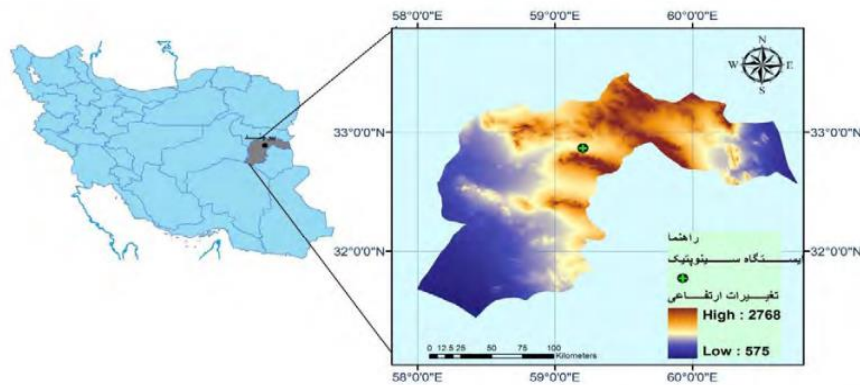
با توجه به مطالعات انجام شده در زمینه‌ی پیش‌بینی بارش، می‌توان به اهمیت برآورد این متغیر در میان سایر پارامترهای هواشناسی پی برد. در این مطالعه پارامترهای مؤثر بر پیش‌بینی بارش ماهانه تعیین و دو مدل ناپارامتریک درخت تصمیم M5 و الگوریتم KNN در راستای برآورد این پدیده‌ی هواشناسی با یکدیگر مقایسه و مدل برتر ارائه گردید، که تاکنون در ایران مقایسه‌ای بین این دو مدل داده کاوی جهت پیش‌بینی بارش ماهانه صورت نگرفته است.

مواد و روش‌ها

ایستگاه سینوپتیک بیرجند با ارتفاعی معادل ۱۴۹۱ متر از سطح دریا در موقعیت جغرافیایی ۱۲° ۵۹' طول شرقی و ۳۲° ۵۲' واقع شده است (شکل ۱).

دسته بندی، شناخت الگوها، پیش‌بینی و مدل‌سازی داده‌ها وجود دارد که در یک نگاه کلی می‌توان این روش‌ها را به دو دسته پارامتری و ناپارامتری تقسیم‌بندی نمود. در مدل‌های ناپارامتری، مرحله تخمین پارامترها وجود ندارد. یکی از شناخت شده‌ترین مدل‌های ناپارامتری روش K نزدیک‌ترین همسایه است. این روش به‌طور گسترده‌ای در علوم مختلف از جمله شناخت الگوهای داده‌ها و دسته‌بندی اطلاعات مورد استفاده قرار گرفته است (Lall and Yakowitz., 1985; Sharma., 1996). همچنین روش‌های داده کاوی برای مجموعه داده‌های بزرگ با متغیرهای زیاد ساخته شده‌اند، بنابراین متفاوت از روش‌های آماری قدیمی هستند که برای مجموعه داده‌های کوچک با متغیرهای اندک طراحی شده‌اند. روش‌های بر مبنای درخت یکی از تکنیک‌های داده کاوی است که در این روش‌ها خروجی به صورت یک مدل با سازه درختی با استفاده از داده‌های ورودی و خروجی می‌باشد. الگوریتم M5 رایج‌ترین طبقه‌بندی استفاده شده در خانواده مدل تصمیم‌گیری درختی است (ستاری و نهرین، ۱۳۹۲).

در سال‌های اخیر مطالعات بسیاری در زمینه‌ی پیش‌بینی بارش و تعیین پارامترهای مؤثر بر آن با استفاده از روش‌های مختلف انجام گردیده است که از این میان می‌توان به مطالعات (Chuan (1997) ، (2002) ، (Trafalis et al., (2005) ، (Maria et al., (2007) ، (2008) و (Hung et al., (2008) اشاره نمود که با استفاده از روش شبکه عصبی مصنوعی مقادیر بارش را برآورد نمودند. در ایران نیز خلیلی و همکاران (۱۳۸۹)، با استفاده از روش شبکه عصبی مصنوعی مقدار بارش ماهانه‌ی ایستگاه سینوپتیک مشهد را پیش‌بینی نمودند که نتایج به دست آمده بر کارایی این روش دلالت داشت. فدایی کرمانی و همکاران (۱۳۹۳)، با استفاده از روش الگوریتم K نزدیکترین همسایگی خشکسالی بر مبنای شاخص بارش استاندارد را در شهرستان بم پیش‌بینی نمودند. نتایج به دست آمده بیان می‌کند که مدل KNN می‌تواند پیش‌بینی‌های قابل قبولی از وضعیت خشکسالی منطقه را ارائه نماید. سیدکابلی و همکاران (۱۳۹۱)، براساس روش ناپارامتریک نزدیکترین همسایگی مدلی جهت ریزمقیاس نمایی داده‌های اقلیمی ارائه نمودند. نتایج به دست آمده نشان داد که روش جمعی مبتنی بر الگوریتم نزدیکترین همسایگی عملکرد بهتری در مقایسه با روش جمعی غیرخطی احتمالاتی دارد و از عدم قطعیت کمتری در پیش‌بینی برخوردار است. قربانی (۱۳۹۳)، کارایی مدل‌های داده‌کاوی ماشین بردار پشتیبان، درخت تصمیم و نزدیکترین همسایگی k را در ریزمقیاس نمایی بارش بر اساس داده‌های مدل گردش عمومی NCEP را بررسی نمود. نتایج این بررسی نشان داد که بارش شبیه‌سازی شده با هر یک از مدل‌های داده‌کاوی، دارای میانگین و انحراف معیار کمتری نسبت به داده‌های



شکل ۱- موقعیت جغرافیایی ایستگاه سینوپتیک بیرجند (هادی زاده و همکاران، ۱۳۹۰)

امین نتیجه‌ی تست پتانسیلی را دارند، Sd بیانگر انحراف معیار، y_i مقدار عددی ویژگی هدف نمونه‌ی i و N شماره‌ی داده‌ها را نشان می‌دهد (Alberg et al., 2012؛ ستاری و نهرین، ۱۳۹۲).

نزدیک‌ترین همسایگی k (KNN)^۲

در حالت کلی از این الگوریتم به دو منظور استفاده می‌شود: برای تخمین تابع چگالی توزیع داده‌های تعلیم و برای طبقه‌بندی داده‌های تست براساس الگوهای تعلیم. بر خلاف توابع انتقالی کلاسیک، مدل نزدیک‌ترین همسایگی از هیچ تابع ریاضی از پیش تعریف شده‌ای برای تخمین متغیرهای مختلف استفاده نمی‌کند. به طور کلی می‌توان گفت که مدل نزدیک‌ترین همسایگی، یکی از روش‌های داده‌کاوی است که هدف کلی آن طبقه‌بندی و تخمین ویژگی‌های یک سری داده‌های مجهول با توجه به بیش‌ترین شباهت این داده‌ها با داده‌های معلومی است که در همسایگی (نزدیکی) آن‌ها قرار دارند (Xindung and Kumar., 2009؛ فدایی کرمانی و همکاران، ۱۳۹۳). در این مدل داده‌های هدف مورد جستجو با توجه به نزدیک‌ترین فاصله نسبت به داده‌های آموزش (بانک داده مرجع)، طبقه‌بندی می‌شوند. نخستین گام در استفاده از این مدل، یافتن روش و رابطه‌ای برای محاسبه‌ی فاصله‌ی بین داده‌های مورد آزمایش و داده‌های تعلیم است. معمولاً برای تعیین این فاصله از فاصله‌ی اقلیدسی زیر استفاده می‌شود (Jagtap et al., 2004):

$$d(X, Y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2} \quad (3)$$

که در آن X نماینده‌ی داده‌های تعلیم با پارامترهای مشخص (x_1) تا x_n و Y نماینده‌ی داده‌های آموزش با همان تعداد پارامترهای مشخص (y_1 تا y_n) است.

درخت تصمیم‌گیری M5^۱

روش درخت تصمیم یک روش سلسله‌مراتبی یا چند مرحله‌ای است که در آن به صورت بازگشتی مجموعه داده‌ها به روش دودویی به تقسیمات فرعی و کوچک‌تر تقسیم‌بندی می‌شود تا زمانی که تقسیمات فرعی نهایی نتوانند بیشتر از آن تجزیه شوند. درختان تصمیم استقرایی مجموعه‌ای از داده‌های معلوم را می‌گیرد و یک درخت تصمیم را از آن استخراج می‌کند. سپس درخت می‌تواند به صورت مجموعه قوانینی برای پیش‌بینی نتیجه ویژگی‌های معلوم استفاده شود (طالبی و اکبری، ۱۳۹۲). اولین مرحله برای ایجاد یک مدل درختی، استفاده از یک معیار انشعاب است. معیار انشعاب برای الگوریتم M5 براساس عملکرد انحراف استاندارد مقادیر هر کلاس و یا طبقه است که در هر گره به دست آمده است. این روش اساس روش‌های طبقه‌بندی است که آنتروپی نامیده می‌شود. آنتروپی می‌تواند به عنوان معیار میزان آشفتگی و بی‌نظمی یک سیستم تفسیر شود. معیار انشعاب بیانگر میزان خطا در آن گره می‌باشد و مدل حداقل خطای مورد انتظار را به عنوان نتیجه‌ی آزمایش هر صفت در آن گره محاسبه می‌کند. خطای مدل عموماً با اندازه‌گیری دقت پیش‌بینی مقادیر هدف موارد دیده نشده سنجش می‌شود. فرمول محاسبه‌ی کاهش انحراف استاندارد (SDR) به صورت روابط زیر می‌باشد:

$$SDR = Sd(T) - \sum_{i=1}^N \frac{|T_i|}{|T|} Sd(T_i) \quad (1)$$

$$Sd(T) = \sqrt{\frac{1}{N} \left(\sum_{i=1}^N y_i^2 - \frac{1}{N} \left(\sum_{i=1}^N y_i \right)^2 \right)} \quad (2)$$

در این رابطه T مجموعه‌ای از نمونه‌هاست (موارد) که به هر گره وارد می‌شود، T_i نشان‌دهنده‌ی زیرمجموعه‌ای از نمونه‌هاست که i

تجسم و تحلیل بسیاری از الگوریتم‌های داده کاوی را دارد (Witten and Frank., 2000).

در این مطالعه ابتدا پارامترهای اقلیمی مؤثر بر بارش مانند اختلاف میانگین حداکثر و حداقل دما ($\bar{T}_{max}-\bar{T}_{min}$)، دمای حداکثر (T_{max}) و حداقل (T_{min})، میانگین رطوبت نسبی ($R\bar{H}$)، میانگین حداکثر ($R\bar{H}_{max}$) و حداقل رطوبت نسبی ($R\bar{H}_{min}$)، میانگین فشار بخار آب ($W\bar{P}$)، میانگین فشار هوا ($A\bar{P}$)، میانگین سرعت باد ($W\bar{S}$)، درجه روز سرمایش (پایه ۲۱ درجه سانتی‌گراد) ($CDD21$)، تعداد روزهای ابری (NCD) و طوفانی (NSD) برای یک دوره‌ی آماری ۵۰ ساله طی سال‌های ۲۰۱۰-۱۹۶۱ میلادی به صورت ماهانه به عنوان داده‌های ورودی به نرم افزار Weka و داده‌های بارش ماهانه به عنوان خروجی نرم‌افزار، از سامانه‌ی اطلاعاتی ایستگاه سینوپتیک بیرجند استخراج گردید. سپس از تعداد ۶۰۰ داده‌ی موجود در هر پارامتر، ۷۰ درصد از داده‌ها معادل ۴۲۰ داده به عنوان داده‌های تولید و ۳۰ درصد از داده‌ها برابر با ۱۸۰ داده به عنوان داده‌های آزمون، در سه حالت داده‌های خام، میانگین متحرک سه ساله و میانگین متحرک پنج ساله داده‌ها در قالب بیست سناریوی تعریف شده در تحقیق (جدول ۱) جهت پیش‌بینی بارش در همان ماه به هر یک از مدل‌های درخت تصمیم M5 و نزدیک‌ترین همسایگی K وارد گردید. پس از اجرای مدل‌ها، سناریوهای فوق در حالت‌های مختلف با استفاده از معیارهای ارزیابی مورد بررسی قرار گرفتند.

$$X = (x_1, x_2, \dots, x_n) \quad (4)$$

$$Y = (y_1, y_2, \dots, y_n) \quad (5)$$

پس از تعیین فاصله‌ی اقلیدسی بین داده‌ها، نمونه‌های بانک داده به ترتیب صعودی از کم‌ترین فاصله (حداکثر تشابه) تا بیشترین فاصله (حداقل تشابه) از نمونه‌ی مورد نظر، طبقه‌بندی و ارزش‌گذاری می‌شوند. قدم بعدی در این مدل یافتن تعداد نقاطی (k) از داده‌های آزمایش برای تخمین ویژگی‌های بانک داده‌ی مورد نظر است. تعیین تعداد همسایه‌ها (k)، یکی از کلیدی‌ترین و مهم‌ترین مراحل در مدل یاد شده به شمار می‌آید و میزان کارایی این روش به طور قابل ملاحظه‌ای به کیفیت انتخاب نزدیک‌ترین (مشابه‌ترین) نمونه‌ها از داده‌های بانک مرجع بستگی دارد (Xindung and Kumar., 2009: فدایی کرمانی و همکاران، ۱۳۹۳).

نرم‌افزار Weka

این نرم افزار اولین بار در سال ۱۹۹۲ به منظور جمع آوری و یکپارچگی الگوریتم‌های یادگیری ماشین و ابزاری برای پردازش داده‌ای با استفاده از JAVA پیاده سازی و به صورت کد باز تحت مجوز عمومی GNU انتشار گردید و در سال ۱۹۹۳ با بکارگیری اهداف داده کاوی ارتقا داده شد. در حال حاضر این نرم افزار حاوی تعداد زیادی از تکنیک‌های یادگیری ماشین و داده کاوی است که امکان مقایسه‌ی تکنیک‌های مختلف یادگیری ماشین را می‌دهد. در این واسط گرافیکی کاربر اجازه‌ی دسترسی به قابلیت‌هایی مانند

جدول ۱- سناریوهای تعریف شده در مدل‌های داده کاوی*

سناریو	پارامترهای ورودی تعریف شده در هر سناریو
1	$\bar{T}_{max}-\bar{T}_{min}, T_{max}, T_{min}, R\bar{H}, R\bar{H}_{max}, R\bar{H}_{min}, W\bar{P}, A\bar{P}, W\bar{S}, CDD21, NCD, NSD$
2	$R\bar{H}, CDD21$
3	$R\bar{H}_{max}, CDD21$
4	$T_{min}, R\bar{H}_{max}, NCD$
5	$R\bar{H}, W\bar{P}, A\bar{P}, W\bar{S}$
6	$T_{min}, CDD21, NSD$
7	$\bar{T}_{max}-\bar{T}_{min}, W\bar{S}, NCD$
8	$T_{max}, R\bar{H}_{max}$
9	$T_{min}, R\bar{H}_{max}, W\bar{P}, CDD21, NCD$
10	$W\bar{P}, A\bar{P}, NCD$
11	$T_{max}, R\bar{H}_{min}, A\bar{P}, W\bar{S}$
12	$T_{min}, R\bar{H}_{max}, W\bar{P}, W\bar{S}, CDD21, NCD, NSD$
13	$\bar{T}_{max}-\bar{T}_{min}, A\bar{P}, W\bar{S}, NCD$
14	$\bar{T}_{max}-\bar{T}_{min}, R\bar{H}, W\bar{P}, A\bar{P}, W\bar{S}, CDD21$
15	$\bar{T}_{max}-\bar{T}_{min}, R\bar{H}, W\bar{S}, CDD21$
16	$\bar{T}_{max}-\bar{T}_{min}, R\bar{H}, W\bar{P}, A\bar{P}, W\bar{S}, CDD21, NCD, NSD$
17	$R\bar{H}, W\bar{S}, CDD21$
18	$T_{min}, R\bar{H}, CDD21$
19	$T_{min}, R\bar{H}_{max}, A\bar{P}, CDD21$
20	$R\bar{H}_{min}, W\bar{P}, A\bar{P}, W\bar{S}, NSD$

* در کلیه‌ی سناریوهای مورد نظر بارش ماهانه به عنوان پارامتر خروجی به مدل معرفی شده است.

معیارهای ارزیابی

جهت ارزیابی مدل‌ها در سه حالت مختلف داده‌های خام، میانگین متحرک سه ساله و میانگین متحرک پنج ساله و همچنین تعیین بهترین سناریو از سه معیار ضریب همبستگی (r)، جذر میانگین مربعات خطا (RMSE) و میانگین قدر مطلق خطا (MAE) استفاده گردید. همچنین برای ارزیابی بهترین سناریو در حالت‌های مورد نظر علاوه بر سه مدل ارزیابی فوق معیار ضریب نش-ساتکلیف (NS) محاسبه گردید.

$$r = \frac{\sum_{i=1}^n (O_i - \bar{O})(P_i - \bar{P})}{\sqrt{\sum_{i=1}^n (O_i - \bar{O})^2} \sqrt{\sum_{i=1}^n (P_i - \bar{P})^2}} \quad (6)$$

$$RMSE = \sqrt{\frac{\sum_{i=1}^N (P_i - O_i)^2}{N}} \quad (7)$$

$$MAE = \frac{1}{N} \sum_{i=1}^N |P_i - O_i| \quad (8)$$

$$NS = 1 - \frac{\sum_{i=1}^N (O_i - P_i)^2}{\sum_{i=1}^N (O_i - \bar{O})^2} \quad (9)$$

در روابط فوق $O_i, \bar{O}_i, P_i, \bar{P}_i$ و N به ترتیب مقادیر مشاهداتی، میانگین مقادیر مشاهداتی، مقادیر پیش‌بینی شده، میانگین مقادیر پیش‌بینی شده و تعداد مشاهدات را نشان می‌دهد. پایین بودن مقادیر خطا (RMSE و MAE) و بالا بودن مقادیر ضریب همبستگی و ضریب نش-ساتکلیف نشان دهنده‌ی دقت مدل و حالت در نظر گرفته شده است.

نتایج و بحث

نتایج ارزیابی سناریوها با استفاده از معیارهای ارزیابی (جدول ۲)، نشان می‌دهد که از میان بیست سناریوی تعریف شده در تحقیق حاضر، سناریوی پانزدهم با متغیرهای ورودی اختلاف میانگین حداکثر و حداقل دما، متوسط رطوبت نسبی، میانگین سرعت باد و درجه روز سرمایه‌ش (بر پایه ۲۱ درجه سانتی‌گراد) ماهانه با مقادیر $R=0/72$ ، $RMSE=13/33$ و $MAE=7/008$ در مدل KNN و $R=0/86$ ،

$RMSE=9/71$ و $MAE=5/57$ در مدل M5 در داده‌های خام، $R=0/87$ ، $RMSE=7/214$ و $MAE=4/655$ در مدل KNN و $R=0/933$ ، $RMSE=5/256$ و $MAE=3/545$ در مدل M5 در داده‌های میانگین متحرک سه ساله و $R=0/888$ ، $RMSE=5/582$ و $MAE=4/112$ در مدل KNN و $R=0/944$ ، $RMSE=3/93$ و $MAE=3/157$ در مدل M5 در داده‌های میانگین متحرک پنج ساله پاسخ بهتری به پیش‌بینی مقادیر بارش در همان ماه می‌دهد. همچنین مقایسه‌ی حالت‌های بیان شده در هر یک از مدل‌ها حاکی از آن است که در تمامی سناریوها به جز دو سناریوی یازدهم و بیستم میانگین متحرک پنج ساله نسبت به میانگین متحرک سه ساله و داده‌های خام پیش‌بینی دقیق‌تری را ارائه می‌دهد. در دو سناریوی یازدهم و بیستم نیز پیش‌بینی بر اساس داده‌های میانگین متحرک سه ساله بهتر از میانگین متحرک پنج ساله و هر دوی این حالات از پیش‌بینی بر اساس داده‌های خام بهتر می‌باشد. بنابراین می‌توان دریافت که پیش‌بینی بر اساس داده‌های میانگین متحرک دقت بیشتری نسبت به داده‌های خام ایستگاه سینوپتیک بیرجند دارد. نتایج بیان شده در زمینه‌ی مقایسه‌ی حالات مختلف داده‌ها و کارایی مدل درختی M5 با نتایج دستورانی و همکاران (۱۳۹۱)، که حاکی از توانایی مدل درخت تصمیم‌گیری خصوصاً در شرایط استفاده از میانگین متحرک ۵ ساله داده‌ها می‌باشد، مطابقت دارد. از طرفی مقایسه‌ی بین دو مدل داده‌کاوی درخت تصمیم و الگوریتم k نزدیک‌ترین همسایه نشان می‌دهد که در تمامی سناریوها در پیش‌بینی بارش ماهانه مدل درختی M5 به طور میانگین با مقادیر $R=0/813$ ، $RMSE=11/588$ و $MAE=7/446$ در داده‌های خام، $R=0/897$ ، $RMSE=7/446$ و $MAE=5/245$ در داده‌های میانگین متحرک سه ساله و $R=0/904$ ، $RMSE=6/054$ و $MAE=4/780$ در داده‌های میانگین متحرک پنج ساله نسبت به الگوریتم KNN به طور میانگین با مقادیر $R=0/719$ ، $RMSE=15/296$ و $MAE=8/493$ در داده‌های خام، $R=0/838$ ، $RMSE=9/071$ و $MAE=5/982$ در داده‌های میانگین متحرک سه ساله و $R=0/837$ ، $RMSE=7/698$ و $MAE=5/595$ در داده‌های میانگین متحرک پنج ساله ارجحیت دارد.

جدول ۲- نتایج ارزیابی مدل‌ها

داده‌های خام	میانگین متحرک ۳ ساله		میانگین متحرک ۵ ساله		
	M5	KNN	M5	KNN	
R	0/84	0/887	0/907	0/876	
سناریو ۱	13/16	7/485	8/669	7/554	RMSE
	7/27	4/796	5/834	5/213	MAE
R	0/72	0/775	0/835	0/817	
سناریو ۲	13/66	12/71	8/465	8/202	RMSE

۵/۵۱۶	۶/۳۴۳	۶/۵۴۸	۶/۹۰۴	۸/۵۶	۷/۳۱	MAE	
-/۸۹	-/۸۴۵	-/۸۶۸	-/۸۰۸	-/۷۶	-/۶۸	R	
۶/۴۴	۷/۷۸۳	۷/۸۷۶	۱۰/۰۶۳	۱۲/۵۲	۱۴/۵۴	RMSE	سناریو ۳
۵/۰۷۵	۵/۸۳۹	۵/۴۶۸	۶/۷۳۹	۷/۳۵	۸/۱۴	MAE	
-/۹۱۸	-/۸۳۳	-/۹۲۳	-/۸۳۹	-/۸۳	-/۶۸	R	
۵/۴۴۵	۷/۴۷۳	۶/۲۷۹	۹/۱۸۷	۱۰/۸	۱۵/۹۴	RMSE	سناریو ۴
۳/۸۷۹	۵/۳۵۵	۴/۳۷۳	۶/۱۰۳	۶/۶۷	۸/۹۳	MAE	
-/۸۹۶	-/۸۲۶	-/۹۰۲	-/۸۶۵	-/۸۴	-/۷۹	R	
۷/۷۱۵	۸/۷۸۲	۹/۶۳۶	۱۰/۷۸۸	۱۱/۵۷	۱۵/۵۷	RMSE	سناریو ۵
۶/۴۴۳	۶/۵۱۳	۶/۹۸	۷/۱۷۸	۷/۳۶	۸/۲۳	MAE	
-/۸۳۹	-/۷۵۱	-/۸۲	-/۷۶۹	-/۷۳	-/۵۸	R	
۷/۸۳	۹/۴۷۱	۹/۷۰۲	۱۰/۳۴	۱۳/۷۷	۱۹/۲۴	RMSE	سناریو ۶
۶/۱۳۴	۶/۸۹	۶/۷۹۹	۶/۷۱۸۸	۹/۳۸	۱۱/۵۴	MAE	
-/۹۳۳	-/۸۴۲	-/۹۳	-/۷۹۲	-/۸۴	-/۷۱	R	
۴/۵۰۴	۶/۶۷۶	۵/۴۸	۸/۹۶۵	۱۰/۳۶	۱۳/۶۳	RMSE	سناریو ۷
۳/۳۶۹	۴/۹۲۵	۳/۹۱۷	۶/۱۵۷	۶/۶۴	۷/۸۱	MAE	
-/۸۸۵	-/۸۰۲	-/۸۹	-/۸۲۷	-/۸	-/۶۶	R	
۶/۵۵۷	۸/۲۲۳	۷/۷۷۳	۱۰/۲۳	۱۲/۵۴	۱۸/۴۸	RMSE	سناریو ۸
۵/۰۸	۶/۱۷۵	۵/۳۴۳	۶/۶۷۶	۷/۶۵	۱۰/۰۹	MAE	
-/۹۲۱	-/۸۶۶	-/۹۱۵	-/۸۵۸	-/۸۱	-/۷۵	R	
۵/۱۷	۶/۵۲	۶/۵۳۸	۸/۴۶۵	۱۱/۲۷	۱۴/۱۸	RMSE	سناریو ۹
۳/۹۵۳	۴/۷۳	۴/۳۷	۵/۴۴	۷/۴۲	۷/۴۷	MAE	
-/۸۵۳	-/۸۲	-/۸۳۹	-/۷۵	-/۷۲	-/۶۷	R	
۶/۲۳	۸/۰۸۶	۸/۰۱۳	۱۰/۵۷۷	۱۳/۰۷	۱۵/۴۷	RMSE	سناریو ۱۰
۴/۵۴۹	۵/۶۸۷	۵/۴۴۹	۶/۸۰۷	۸/۱۷	۸/۹	MAE	
-/۸۶۶	-/۸۲۱	-/۸۷۷	-/۸۸۱	-/۸۱	-/۷۶	R	
۸/۷۳۲	۸/۶۱۱	۹/۵۱۷	۸/۹۰۳	۱۳/۱۷	۱۶/۹۹۶	RMSE	سناریو ۱۱
۷/۳۸۸	۶/۴۷۵	۷/۱۷۲	۵/۹۲۲	۹/۰۵	۹/۷۵	MAE	
-/۹۲۲	-/۸۷	-/۹۱۷	-/۸۷۱	-/۸۵	-/۷۲	R	
۵/۲۲۱	۶/۰۹	۶/۸۳۳	۷/۸۷۴	۱۰/۶۵	۱۴/۶۳	RMSE	سناریو ۱۲
۳/۷۹۲	۴/۲۷۱	۴/۴۱۸	۵/۰۹	۶/۶۷	۸/۱۱	MAE	
-/۹۳۳	-/۸۶۱۸	-/۹۳۳	-/۸۶۶	-/۷۸	-/۷۱	R	
۴/۴۲۵	۶/۴۷۸	۶/۰۰۷	۷/۸۹۱	۱۱/۷۸	۱۶/۶۹	RMSE	سناریو ۱۳
۳/۲۹۶	۴/۵۴۵	۴/۰۵۵	۵/۳۱۷	۸/۱۴	۹/۴۴	MAE	
-/۹۴۴	-/۸۶۹	-/۹۳۴	-/۹۰۸	-/۸۶	-/۸	R	
۴/۹۹۲	۷/۵۱	۶/۱۴	۷/۹۴۲	۱۱/۴۴	۱۳/۱۶	RMSE	سناریو ۱۴
۳/۹۶	۵/۴۵۴	۴/۳۳	۵/۲۱۹	۷/۱۷	۷/۳۷	MAE	
-/۹۴۴	-/۸۸۸	-/۹۳۳	-/۸۷	-/۸۶	-/۷۲	R	
۳/۹۳	۵/۵۸۲	۵/۲۵۶	۷/۲۱۴	۹/۷۱	۱۳/۳۳	RMSE	سناریو ۱۵
۳/۱۵۷	۴/۱۱۲	۳/۵۴۵	۴/۶۵۵	۵/۵۷	۷/۰۰۸	MAE	
-/۹۳۶	-/۸۷۴	-/۹۳۸	-/۸۸۶	-/۸۷	-/۷۶	R	
۵/۶۲۷	۷/۳۹۵	۶/۲۸۸	۷/۶۸۹	۱۰/۱۹	۱۴/۱۳	RMSE	سناریو ۱۶
۴/۱۸۷	۵/۱۰۸	۴/۲۶۸	۵/۰۲۷	۶/۵۲	۷/۶۳	MAE	
-/۸۷۸	-/۷۹۲	-/۸۵۶	-/۷۵۶	-/۸	-/۷۴	R	
۵/۷۵۴	۸/۱۱	۷/۶۰۸	۱۰/۱۹۷	۱۱/۳۷	۱۲/۹۲	RMSE	سناریو ۱۷
۴/۵۸۹	۵/۸۸۶	۵/۷۵	۶/۴۹۷	۶/۵۴	۶/۷۵	MAE	

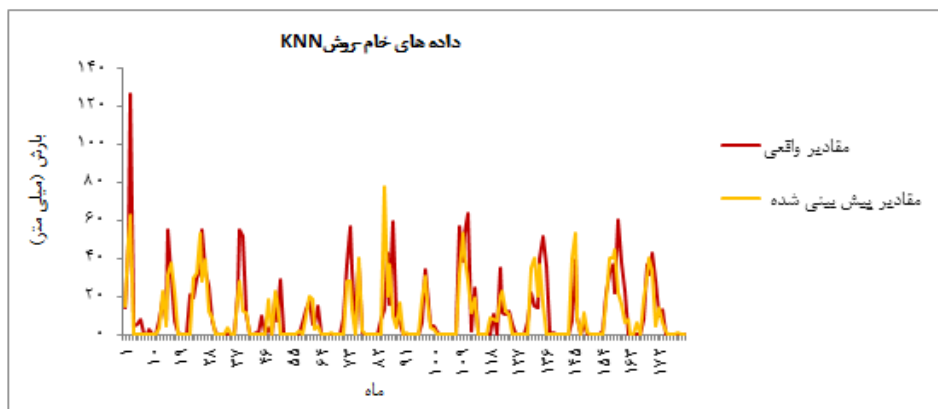
۰/۸۹۳	۰/۷۹۵	۰/۹	۰/۸۲۹	۰/۸۳	۰/۶۷	R	
۶/۲۰۹	۸/۸۴۵	۷/۱۱۳	۹/۰۶	۱۱/۲۵	۱۶/۸۵	RMSE	سناریو ۱۸
۴/۹۶۷	۶/۵۴۲	۴/۹۷۱	۶/۱۷۲	۶/۸۱	۹/۱۲	MAE	
۰/۹۱۸	۰/۸۵۲	۰/۹۱۲	۰/۸۶۴	۰/۸۴۵	۰/۷۳	R	
۶/۴۷۷	۷/۹۸۱	۷/۳۹۸	۹/۱۹۷	۱۱/۵۴	۱۷/۴۵	RMSE	سناریو ۱۹
۵/۱۹۲	۵/۷۶۵	۵/۲۲۶	۶/۱	۷/۱۸	۱۰/۱۲	MAE	
۰/۹۲۳	۰/۸۳۷	۰/۹۲۴	۰/۸۶۷	۰/۸۴	۰/۷۶	R	
۷/۶۴۷	۸/۵۹۳	۸/۳۰۴	۹/۲۶۱	۱۱/۴	۱۵/۸۹	RMSE	سناریو ۲۰
۶/۶۵۹	۶/۰۷۲	۶/۰۷۷	۶/۱۲۱	۸/۵	۸/۸۸	MAE	

می‌توان دریافت که در مدل درختی M5 روند مقادیر پیش‌بینی شده تطابق بیشتری بر روند داده‌های واقعی بارش ماهانه نسبت به روش KNN دارد. از طرفی مقایسه‌ی نمودارها در حالت‌های مختلف داده‌های خام، میانگین متحرک سه ساله و میانگین متحرک پنج ساله، نشان می‌دهد که میانگین متحرک پنج ساله به علت نوسان کمتری که در پیش‌بینی بارش دارد، نسبت به سایر حالات بهترین عملکرد را در پیش‌بینی این پارامتر اقلیمی داراست. (شکل‌های ۲-۷)

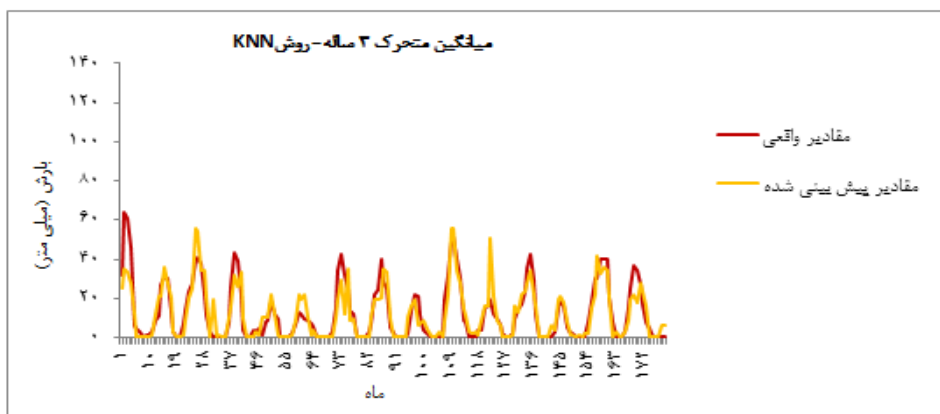
نتایج به دست آمده از برآورد ضریب نش-ساتکلیف برای مقایسه‌ی مدل‌ها و حالات مختلف داده‌ها در سناریوی برتر (جدول ۳) نیز بر دقت بیشتر مدل M5 نسبت به الگوریتم KNN دلالت دارد و همچنین نشان می‌دهد که در بین حالات مختلف داده‌ها پیش‌بینی بارش ماهانه با استفاده از داده‌های میانگین متحرک پنج ساله در مقایسه با دو حالت دیگر بیشتر به واقعیت نزدیک است. با مقایسه‌ی متناظر نمودار مدل‌های مختلف در هر یک از حالات داده‌های ورودی،

جدول ۳- نتایج ارزیابی سناریو ۱۵ با ضریب Nash-Sutcliffe

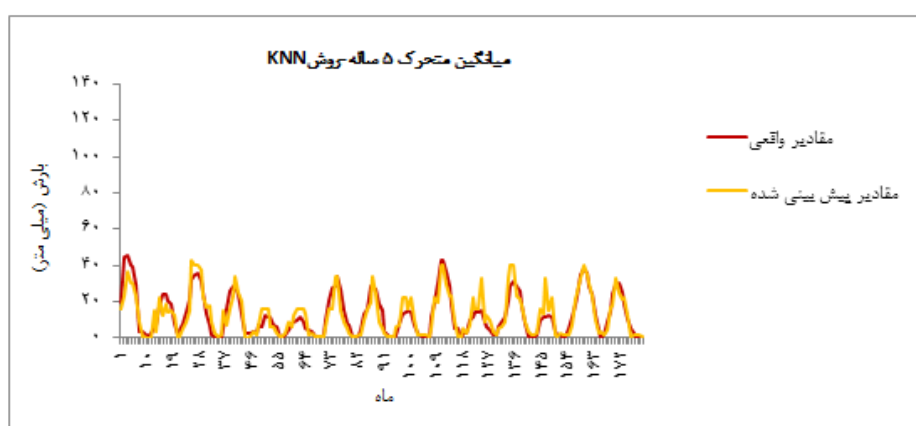
میانگین متحرک ۵ ساله		میانگین متحرک ۳ ساله		خم		نوع داده‌های ورودی
M5	KNN	M5	KNN	M5	KNN	مدل
۰/۸۹۳	۰/۷۷۷	۰/۸۷۷	۰/۷۵۳	۰/۷۴۳	۰/۴۹۸	ضریب نش-ساتکلیف سناریو ۱۵



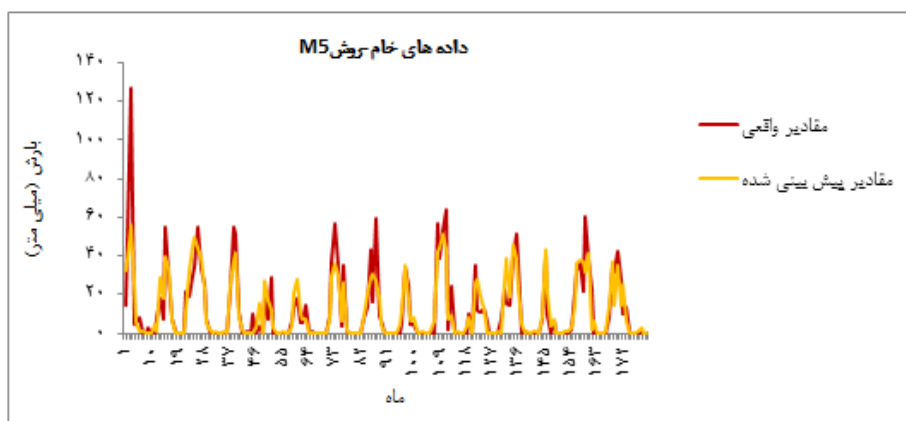
شکل ۲- مقایسه‌ی مقادیر واقعی (داده‌های خام) و پیش‌بینی شده‌ی بارش ماهانه با روش KNN- سناریو ۱۵



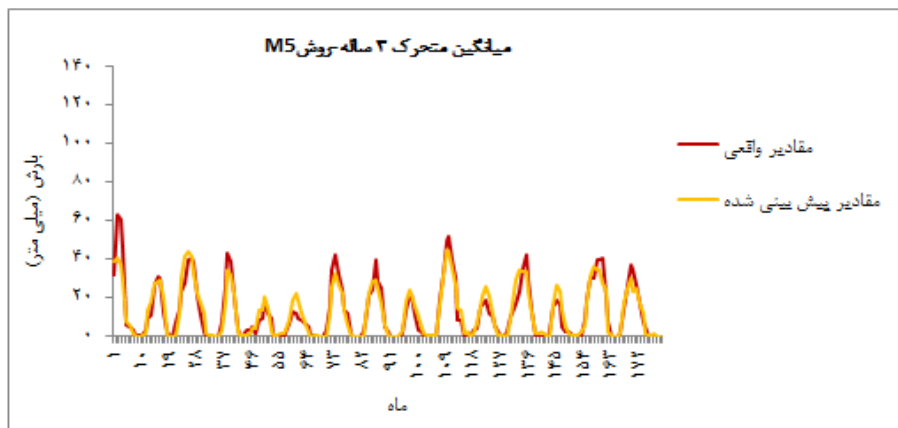
شکل ۳- مقایسه‌ی مقادیر واقعی (میانگین متحرک ۳ ساله) و پیش‌بینی شده‌ی بارش ماهانه با روش KNN- سناریو ۱۵



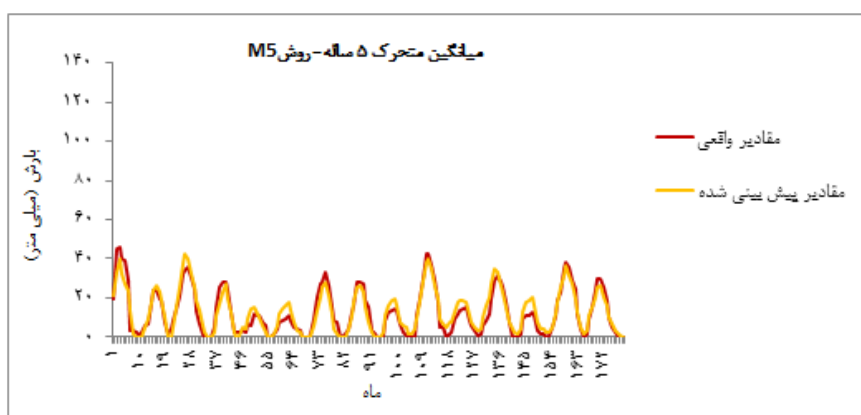
شکل ۴- مقایسه‌ی مقادیر واقعی (میانگین متحرک ۵ ساله) و پیش‌بینی شده‌ی بارش ماهانه با روش KNN- سناریو ۱۵



شکل ۵- مقایسه‌ی مقادیر واقعی (داده‌های خام) و پیش‌بینی شده‌ی بارش ماهانه با روش M5- سناریو ۱۵



شکل ۶- مقایسه‌ی مقادیر واقعی (میانگین متحرک ۳ ساله) و پیش‌بینی شده‌ی بارش ماهانه با روش M5- سناریو ۱۵



شکل ۷- مقایسه‌ی مقادیر واقعی (میانگین متحرک ۵ ساله) و پیش‌بینی شده‌ی بارش ماهانه با روش M5- سناریو ۱۵

نتیجه‌گیری

امروزه با توجه به خسارات قابل توجهی که در اثر وقوع پدیده‌های هیدرولوژیکی مانند خشکسالی و سیل به منابع طبیعی و انسانی وارد می‌گردد، سازماندهی و حفاظت از این منابع، نیازمند برنامه‌ریزی و تعیین موقعیت زمانی و مکانی وقوع این پدیده‌هاست، که پیش‌بینی بارندگی به عنوان یکی از مهم‌ترین عوامل مؤثر در برآورد این پدیده‌ها مطرح می‌گردد. این مطالعه با هدف پیش‌بینی بارش ماهانه از طریق ارائه‌ی یک مدل داده‌کاوی کارا و همچنین تعیین پارامترهای مؤثر در افزایش دقت در پیش‌بینی این پارامتر هواشناسی انجام گردیده است. به طور کلی نتایج تحقیق حاضر بیانگر این مطلب است که مدل درختی M5 به عنوان یک مدل داده‌کاوی می‌تواند بارش ماهانه را با دقت بیشتری نسبت به الگوریتم KNN برآورد نماید که این مهم در مواقع استفاده از میانگین متحرک داده‌ها با دقت بیشتری نسبت به داده‌های خام قابل دستیابی است. از طرفی با توجه به ارزیابی سناریوهای مختلف توصیه می‌گردد که از ترکیب متغیرهای ورودی اختلاف میانگین حداکثر و حداقل دما، متوسط رطوبت نسبی، میانگین سرعت باد و درجه روز سرماش (بر پایه ۲۱ درجه سانتی‌گراد) برای

پیش‌بینی بارش ماهانه، جهت تعیین زمان و مکان تقریبی وقوع پدیده‌هایی چون خشکسالی و سیلاب بر مبنای وقوع بارش، استفاده گردد.

منابع

- امیدوار، ک.، شفیعی، ش.، تقی‌زاده، ز. و علی‌پور، م. ۱۳۹۳. ارزیابی کارایی مدل درخت تصمیم در پیش‌بینی بارش ایستگاه سینوپتیک کرمانشاه. نشریه تحقیقات کاربردی علوم جغرافیایی. سال چهاردهم، شماره ۳۴، صفحات ۸۹-۱۱۰.
- خلیلی، ن.، خداشناس، س.ر.، داوری، ک. و موسوی‌بایگی، م. ۱۳۸۹. پیش‌بینی بارش روزانه با استفاده از شبکه‌های عصبی مصنوعی مطالعه موردی: ایستگاه سینوپتیک مشهد. مجله پژوهش‌های آبخیزداری (پژوهش و سازندگی). شماره ۸۹، صفحات ۷-۱۵.
- دستورانی، م.ت.، حبیبی‌پور، ا.، اختصاصی، م.ر.، طالبی، ع. و محجوبی، ج. ۱۳۹۱. بررسی کارایی مدل درخت تصمیم در پیش‌بینی بارش مطالعه موردی ایستگاه سینوپتیک یزد. مجله تحقیقات منابع آب

69-78.

ایران. سال هشتم، شماره ۳، ۱۴-۲۷.

- Chattopadhyay, S. 2007. Feed forward artificial neural network model to predict the average summer-monsoon rainfall in India, *Acta Geophysical*, No. 55(3), pp. 369-382.
- Chuan, C.S. 1997. Weather prediction using artificial neural network, *Journal of Hydrology*, 230: 101-119.
- Hung, NQ., Babel, MS., Weesakul, S., Tripathi, NK. 2008. An artificial neural network model for rainfall forecasting in Bangkok, Thailand, *Hydrology and Earth System Sciences Discussions*, No. 5, pp. 183-218.
- Jagtap SS, Lall U, Jones, JW, Gijsman AJ, Ritchie JT. 2004. Dynamic nearest-neighbor method for estimating soil water parameters. *Trans. ASAE*. 47:1437-1444.
- Lall, U., and Sharma, A. 1996. A nearest neighbor bootstrap for resampling hydrologic time series, *Water Resources Research*, 32(3), 679-694.
- Maria, C., Haroldo, F., Ferreira, N. 2005. Artificial neural network technique for rainfall forecasting applied to the Sao Paulo region, *Journal of Hydrology*, No. 301, pp.146-162.
- Nash, J. E. and Sutcliffe, J. V. 1970. River flow forecasting through conceptual models, Part I - A discussion of principles, *J. Hydrol.*, 10, 282-290.
- Trafalis, TB., White, A., Santosa, B., Richman, MB. 2002. Data mining techniques for improved WSR-88D rainfall estimation, *Computers in Industrial Engineering*, No. 43, pp. 775-786.
- Witten, L. and Frank, E. 2000. *Data Mining: Practical Machine Learning Tools and Techniques with Java Implementations*, Morgan Kaufmann Publishers.
- Xindung, W. & Kumar, V. 2009. *Top Ten Algorithm in Data Mining*, First Edition, Taylor & Francis Group, USA.
- Yakowitz, S. J. 1985. Nonparametric density estimation, prediction, and regression for markov sequences. *J. Am. Stat. Assoc.*, 80, 215-221.
- ستاری، م.ت. و نهرین، ف. ۱۳۹۲. پیش‌بینی مقادیر حداکثر بارش روزانه با استفاده از سیستم‌های هوشمند و مقایسه آن با مدل درختی M5؛ مطالعه موردی ایستگاه‌های اهر و جلفا. فصلنامه علمی-پژوهشی مهندسی آبیاری و آب. سال چهارم، شماره چهاردهم، صفحات ۸۳-۹۸.
- سیدکابلی، ح.، آخوندعلی، م.ع.، مساح‌بوانی، ح. و رادمش، ف. ۱۳۹۱. ارائه مدل ریزمقیاس‌نمایی داده‌های اقلیمی براساس روش ناپارامتریک نزدیکترین همسایگی (K-NN). نشریه آب و خاک (علوم و صنایع کشاورزی). جلد ۲۶، شماره ۴، صفحات ۸۰۸-۷۷۹.
- طالبی، ع. و اکبری، ز. ۱۳۹۲. بررسی کارایی مدل درختان تصمیم-گیری در برآورد رسوبات معلق رودخانه‌ای (مطالعه موردی: حوضه سد ایلام). *مجله علوم و فنون کشاورزی و منابع طبیعی*، علوم آب و خاک. سال هفدهم، شماره ۶۳، صفحات ۱۲۱-۱۰۹.
- فدایی کرمانی، ا.، خانجانی، م.ج. و بارانی، غ.ع. ۱۳۹۳. کاربرد الگوریتم K-نزدیک‌ترین همسایگی در پایش خشکسالی بر مبنای شاخص بارش استاندارد (SPI) (مطالعه موردی: شهرستان بم). فصلنامه بین‌المللی پژوهشی تحلیلی منابع آب و توسعه. شماره ۱، صفحات ۱۳۸-۱۳۱.
- قربانی، خ. ۱۳۹۳. ارزیابی مدل‌های داده‌کاوی در ریزمقیاس‌نمایی بارش براساس داده‌های مدل‌گردش عمومی NCEP (مطالعه موردی: ایستگاه سینوپتیک کرمانشاه). *مجله پژوهش آب ایران*. سال هشتم، شماره ۱۵، صفحات ۱۷۷-۱۸۶.
- مهدوی، م. ۱۳۷۴. هیدرولوژی کاربردی، چاپ دوم، انتشارات دانشگاه تهران، تهران.
- نعیمی، م. و احقافی، ا. ۱۳۸۱. بررسی و مدیریت خشکسالی در ایران. مرکز اطلاعات و مدارک علمی ایران. شماره ۴۸۴۸۲۴.
- Alberg, D., M. Last and A. Kindle. 2012. Knowledge discovery in data streams with regression tree methods. *WIREs Data Mining Knowl Discov* (2):

Comparison of Decision Tree M5 and K-Nearest Neighborhood Algorithm Models in the Prediction of Monthly Precipitation (Case Study: Birjand Synoptic Station)

F. Poursalehi^{1*}, A. Shahidi², A. Khasheisiuki³
Received: Mar.08, 2019 Accepted: Jun.24, 2019

Abstract

Due to the location of Iran in dry and semi-arid climate, heterogeneous distribution of precipitation and also the occurrence of a climate change phenomenon has caused phenomena such as floods, drought, desertification and dust production and also creating the different economic, social and environmental damages. One of the primary strategies to reduce these losses, is prediction of the precipitation events. The goal of the present study is monthly precipitation prediction with using data mining methods of decision tree (M5) and K-Nearest Neighbor (KNN) algorithms and Comparing these methods in order to determining more efficient method in the field of predicting the precipitation using monthly meteorological data of Birjand synoptic station during the statistical period 1961-2010 in three cases the raw data, the three-year moving average and the five-year moving average in the Weka software. The results showed that in all defined scenarios, the tree model M5 has more ability than the KNN model to predict the monthly precipitation of the station. Also after investigation of the evaluation criteria R, RMSE, MAE and NS, the fifteenth scenario with input variables such as mean difference of maximum and minimum temperature, average relative humidity, average wind speed and cooling degree days (base 21 ° C) in every month was determined as the best scenario for predicting the same month precipitation. Also the obtained results from comparing the defined scenarios in each model in three states raw data, three-year moving average and five-year moving average show that in most scenarios The five-year moving average on average, with the values of R=0.904, RMSE=6.054 and MAE=4.780 in the M5 model and on average, with the values of R=0.837, RMSE=7.698 and MAE=5.595 in the KNN model offers more accurate prediction of monthly Precipitation.

Keywords: Data mining methods, Decision tree, Drouth, K-Nearest Neighborhood, Weka software

1- Ph.D. Student ,Water Resources Engineering at the University of Birjand

2- Associate Professor, Department of Water Science and Engineering University of Birjand

(* - Corresponding Author Email: Fatemehpoursalehi@birjand.ac.ir)