

## مقایسه کارایی مدل‌های درختی در محاسبه ضریب پراکندگی طولی آلاینده‌ها در آبراهه‌های مستقیم

حسین نزارتیان<sup>۱</sup>، جواد ظهیری<sup>۲</sup>، سید محمد کاشفی پور<sup>۳</sup>

تاریخ ارسال ۱۳۹۵/۰۹/۲۱

تاریخ پذیرش ۱۳۹۶/۰۳/۱۳

### چکیده

مدل‌سازی پیشروی آلاینده‌ها در آبراهه‌های طبیعی یکی از مهم‌ترین مسائل محیط زیست است. ضریب پخشیدگی طولی یکی از پارامترهای اساسی در مدل‌سازی انتشار آلودگی‌ها به حساب می‌آید. طی پژوهش‌های صورت گرفته توسط محققان مختلف روابط متعددی جهت برآورد این ضریب ارائه شده است که اغلب این روابط به صورت تجربی و یا نیمه تجربی به دست آمده‌اند. با این وجود، نیاز به روش‌های دقیق‌تر تخمین ضریب پخشیدگی طولی همچنان احساس می‌شود. در این تحقیق جهت تخمین این ضریب، مدل‌های داده‌کاوی با توجه به اطلاعات هیدرولیکی و هندسی رودخانه‌ها توسعه یافته است. بر این اساس الگوریتم‌های درختی CART، M5 و برنامه‌ریزی ژنتیک (GP) مورد استفاده قرار گرفت. جهت مقایسه کارایی مدل‌ها با معادلات موجود از پارامترهای آماری جذر میانگین مربعات خطا، میانگین خطای مطلق و نسبت اختلاف استفاده گردید. نتایج تحلیل‌های آماری نشان داد که مدل‌های داده‌کاوی می‌توانند ضریب پخشیدگی طولی را با دقت بهتر برآورد نمایند. مدل CART با وجود دقت زیاد در مرحله آموزش، در مرحله صحت سنجی از دقت کمتری برخوردار بوده است. مدل‌های M5 و GP به ترتیب دارای جذر میانگین مربعات خطای ۰/۴۱ و ۰/۴۴ و معیار دقت ۶۱٪ و ۶۲٪ بوده و در مقایسه با روابط تجربی موجود از دقت بیشتری برخوردار می‌باشند. با توجه به اختلاف ناچیز میان این دو مدل و سادگی مدل ارائه شده توسط M5، از این مدل می‌توان جهت برآورد ضریب پراکندگی طولی در رودخانه‌ها استفاده کرد.

واژه‌های کلیدی: پیشروی آلاینده‌ها، ضریب پراکندگی طولی، CART، M5، GP

<sup>۱</sup> کارشناسی ارشد، گروه مهندسی آب دانشگاه کشاورزی و منابع طبیعی رامین خوزستان، ۰۹۱۶۹۵۴۸۷۶۸، Hosein.Nezaratian69@gmail.com

<sup>۲</sup> استادیار، گروه مهندسی آب دانشگاه کشاورزی و منابع طبیعی رامین خوزستان، ۰۹۱۶۶۵۳۱۸۹۶، Zahiri\_Javad@yahoo.com (نویسنده مسئول)

<sup>۳</sup> استاد، گروه مهندسی آب دانشگاه شهید چمران اهواز، ۰۹۱۲۱۷۱۵۷۵۲، Kashefipour@excite.com

## سال هشتم • شماره بیست و نهم • پاییز ۱۳۹۶

شدت اختلاط در جهت‌های عرضی و عمقی به تعادل می‌رسد، بعد طولی ضریب پراکندگی بیشتر مورد توجه محققان قرار گرفته است (Chatila, 1997). دامنه تغییرات ضریب انتشار طولی در رودخانه‌های طبیعی با توجه به خصوصیات جریان و هندسه مقطع بسیار متغیر و پیچیده است. اگرچه روابط تجربی و تئوری بسیاری جهت تخمین مقدار ضریب پراکندگی طولی پیشنهاد شده است ولی با این وجود مطالعات بسیاری نیز بر روی این ضریب در حال انجام است (Noori et al., 2009). روابط ارائه شده جهت برآورد ضریب پراکندگی طولی اغلب به صورت تجربی یا نیمه تجربی بوده و در رودخانه‌هایی با شرایط متفاوت دارای دقت‌های متفاوتی هستند. تعدادی از روابط تجربی مورد استفاده جهت برآورد ضریب پخشیدگی طولی در جدول (۱) ارائه شده است. معادلات ارائه شده جهت برآورد ضریب پخشیدگی طولی از پارامترهای مختلف هیدرولیکی و هندسی آبراهه استفاده می‌کنند. در جدول (۱)،  $w$  و  $h$  به ترتیب عرض سطح مقطع جریان و عمق جریان برحسب متر،  $u$  سرعت جریان و  $u^*$  سرعت برشی جریان برحسب متر بر ثانیه و  $K_x$  ضریب پراکندگی طولی برحسب مترمربع بر ثانیه است. در سال‌های اخیر استفاده از الگوریتم‌های هوشمند در شاخه‌های مختلف مهندسی از جمله مهندسی رودخانه رو به افزایش بوده است. الگوریتم‌های فراکاوشی توانایی تخمین پدیده‌ها و فرآیندهای پیچیده طبیعی را دارا بوده و اغلب دارای دقت بیشتری نسبت به روابط تجربی می‌باشند. Etemad-shahidi and Ghaemi (2011) جهت بررسی کارایی الگوریتم M5 در محاسبه عمق آبستگي اطراف مجموعه پایه‌های پل از اطلاعات آزمایشگاهی و میدانی استفاده کردند. Etemad-shahidi and Taghipour (2012) با استفاده از الگوریتم M5 و با به کارگیری اطلاعات ۱۴۹ رودخانه موجود در سرتاسر جهان به تخمین ضریب پراکندگی طولی پرداختند. این محققین جهت تخمین این ضریب از پارامتر سینوسی بودن رودخانه‌ها نیز استفاده کردند. نتیجه این تحقیق بیان گر دقت و قدرت بالای این الگوریتم در تخمین ضریب انتشار طولی است. Sattar and Gharehbaghi (2015) با استفاده از الگوریتم ژنتیک و اطلاعات مربوط به ۱۵۰ رودخانه موجود

## مقدمه

در چند دهه اخیر نحوه عملکرد آبراهه‌های طبیعی و شناسایی رفتار آن‌ها از جهات مختلفی مورد توجه قرار گرفته است که از آن جمله می‌توان به نحوه رفتار آبراهه‌ها در هنگام انتشار آلودگی اشاره کرد. این موضوع به منظور حفظ سلامت عمومی در مناطقی که صنایع بزرگ در کنار رودخانه‌ها قرار گرفته و رودخانه‌ها تامین کننده نیاز اصلی آبی می‌باشند، بیشتر مورد توجه قرار می‌گیرد. بدین جهت توانایی تخمین، شبیه سازی جریان و انتقال آلودگی و رسوب در سیستم مجاری روباز و رودخانه‌ها جهت برنامه‌ریزی منابع آب از اهمیت بالایی برخوردار است (Li and et al., 1998). با تزریق آلاینده‌های مختلف به بالادست یک آبراهه، انتقال آن به پایین دست توسط فرآیندهای اختلاط و انتشار<sup>۱</sup> طولی، عرضی و عمقی رخ می‌دهد (Tayfur and Singh, 2005). اولین بار Taylor (1954) پی برد که سرعت برشی و اختلاط در عرض جریان پس از گذشت زمانی خاص به تعادل رسیده اما این فرآیند در طول مسیر ادامه دارد. او با استفاده از روش فیک برای شار جرمی متلاطم معادله یک بعدی انتقال - پخش را برای کانال‌های یکنواخت ارائه کرد.

$$\frac{\partial C}{\partial t} + u \frac{\partial C}{\partial x} = K_x \frac{\partial^2 C}{\partial x^2} \quad (1)$$

در معادله ۱،  $C$  متوسط غلظت،  $u$  متوسط سرعت طولی،  $t$  زمان،  $x$  جهت طولی در راستای جریان و  $K_x$  ضریب پراکندگی طولی<sup>۲</sup> است. این معادله تمام ویژگی‌های جریان را در کانال‌های آزمایشگاهی و رودخانه‌ها به طور کامل نشان می‌دهد (Ratherford, 1994). ضریب پراکندگی معرفی شده در معادله ۱ نقش مهمی در مدل سازی فرآیند نشت و طراحی آبگیرها داشته و می‌توان آن را نماینده شدت اختلاط در رودخانه‌ها در نظر گرفت (Deng et al., 2002). توانایی جریان رودخانه و دیگر جریان‌های سطحی در پراکنش مواد تزریق شده در جهات طولی، عرضی و عمقی توسط ضرایب پراکندگی  $K_z$ ،  $K_y$ ،  $K_x$  بیان می‌گردد (Tayfur and Singh, 2005). با توجه به این نکته که اندکی پس از انتشار

<sup>1</sup> Dispersion<sup>2</sup> Longitudinal dispersion coefficient

از داده‌های به‌کار رفته در تحقیق Etemad-shahidi and Taghipour (2012) استفاده شده است. تعداد کل این داده‌ها ۱۴۹ سری است که این اطلاعات از رودخانه‌ها و آبراهه‌های سراسر جهان برداشت شده‌اند. این مجموعه شامل اطلاعات هندسی و هیدرولیکی رودخانه‌ها از جمله عرض سطح آب بر حسب متر، عمق جریان بر حسب متر، سرعت جریان و سرعت برشی بر حسب متر بر ثانیه می‌باشند. مشخصات آماری داده‌های مورد استفاده برای آموزش و صحت‌سنجی مدل‌ها در جدول (۲) ارائه شده است.

ضریب پراکندگی طولی از پارامترهای متعدد هندسی و هیدرولیکی آبراهه تأثیر می‌پذیرد که موارد تأثیرگذار بر روی این ضریب را می‌توان به صورت تابع زیر نوشت:

$$K_x = f_1(\rho, \mu, w, h, u, u_*, s_f, s_n) \quad (2)$$

در این رابطه  $\rho$  چگالی مایع،  $\mu$  لزوجت دینامیکی،  $w$  عرض سطح مقطع،  $h$  عمق جریان،  $u$  سرعت جریان،  $u_*$  سرعت برشی جریان،  $s_f$  شیب طولی بستر و  $s_n$  ضریب سینوسی آبراهه است. اغلب پارامترها دارای بعد بوده که جهت جلوگیری از به وجود آمدن ناهمگنی ابعادی در دو طرف معادله از تئوری باکینگهام استفاده گردید (Seo and Cheong, 1998)

$$\frac{K_x}{hu_*} = f_2\left(\rho \frac{uh}{\mu}, \frac{w}{h}, \frac{u}{u_*}, s_f, s_n\right) \quad (3)$$

جریان در رودخانه‌ها و آبراهه‌ها غالباً به صورت متلاطم بوده و می‌توان از عدد رینولدز ( $\rho uh/\mu$ ) صرف نظر نمود. علاوه بر این اندازه‌گیری ضریب سینوسی دشوار و پیچیده بوده و در اغلب رودخانه‌ها این پارامتر اندازه‌گیری نمی‌شود. بر این اساس می‌توان عوامل مؤثر بی‌بعد بر ضریب پراکندگی طولی را به شکل زیر نوشت (Seo and Cheong, 1998; Seo and Beak, 2004)

$$\frac{K_x}{hu_*} = f_2\left(\frac{w}{h}, \frac{u}{u_*}\right) \quad (4)$$

در کشورهای آمریکا، کانادا، اروپا و نیوزلند به تخمین ضریب پراکندگی طولی پرداختند که نتیجه آن گویای برتری و دقت بالای آن نسبت به مدل‌های تجربی بوده است. ظهیری (۱۳۹۴) جهت تخمین عمق آبشستگی اطراف پایه‌های پل از دو الگوریتم M5 و CART استفاده کرد که نتیجه آن بیانگر کارایی بهتر و ساده‌تر الگوریتم M5 بوده است. (Haghiabi (2016) به پیش‌بینی ضریب پراکندگی طولی با استفاده از رگرسیون چند متغیره اسپلاین<sup>۱</sup> (MARS) پرداخت که نتایج به‌دست آمده نشان‌دهنده توانایی مدل اسپلاین در پیش‌بینی این ضریب می‌باشد. هدف از تحقیق پیش رو مقایسه سه الگوریتم درختی M5، GP و CART در پیش‌بینی ضریب پراکندگی طولی و مقایسه نتایج مدل‌های ارائه شده با روابط تجربی است.

جدول (۱): روابط تجربی ارائه شده توسط محققین

محقق	رابطه ارائه شده
Elder (1959)	$K_x = 5.93hu_*$
Fischer (1967)	$K_x = 0.011(u^2w^2/hu_*)$
McQuivey and Keefer (1974)	$K_x = 0.58(h/u_*)^2uw$
Liu (1977)	$K_x = 0.18(u/u_*)^{0.5}(w/h)^2hu_*$
Iwasa and Aya (1991)	$K_x = 2(w/h)^{1.5}hu_*$
Li et al. (1998)	$K_x = 0.55(wu_*/h^2)$
Seo and Cheong (1998)	$K_x = 5.92(u/u_*)^{1.43}(w/h)^{0.62}hu_*$
Koussis and Rodriguez-Mirasol (1998)	$K_x = 0.6(w/h)^2hu_*$
Li et al. (1998)	$K_x = 5.92(u/u_*)^{1.2}(w/h)^{1.3}hu_*$
Kashefipour and Falconer (2002)	$K_x = 10.612(u/u_*)hu_*$
Rajeev and Dutta (2009)	$K_x = 2(u/u_*)^{0.96}(w/h)^{1.25}hu_*$

## مواد و روش‌ها

در تحقیق حاضر جهت مدل‌سازی ضریب پراکندگی طولی

<sup>1</sup> Multivariate Adaptive Regression Splines

جدول (۲): اطلاعات مورد استفاده جهت آموزش و صحت سنجی مدل‌ها

متغیر آماری	عرض (متر)	عمق (متر)	سرعت (متر بر ثانیه)	سرعت برشی (متر بر ثانیه)	ضریب انتشار طولی (مترمربع بر ثانیه)
حداکثر	۲۵۳/۶	۸/۲	۱/۷۳	۰/۵۵	۱۴۸۶/۵
حداقل	۱/۴	۰/۱۴	۰/۰۲۹	۰/۰۰۱۶	۰/۲
میانگین	۴۹/۵۸	۱/۳۵	۰/۴۷	۰/۰۸۴	۸۳/۲۹
انحراف معیار	۴۸/۲۸	۱/۳۲	۰/۳۱	۰/۰۷۲	۱۸۰/۹۵

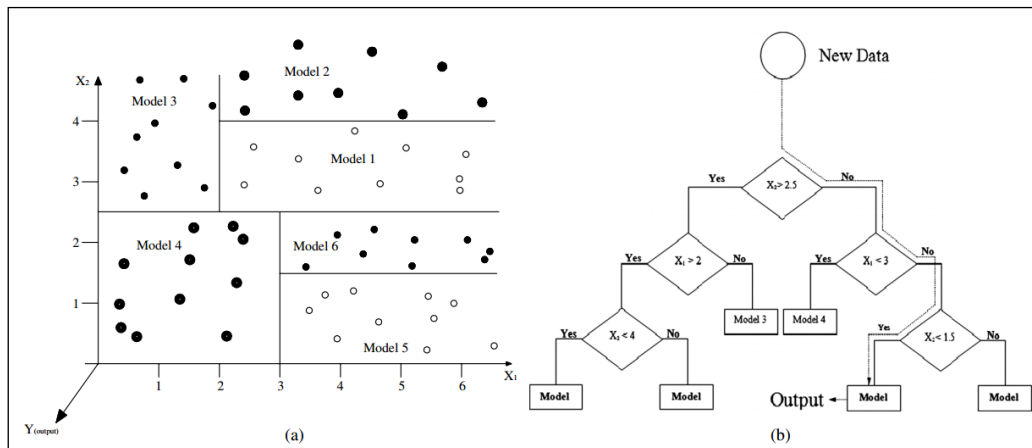
الگوریتم M5 تمامی حالت‌های مختلف جهت ایجاد شاخه بر اساس صفت خاص را بررسی کرده و در نهایت حالتی را انتخاب می‌کند که بتواند بیشتر از حالت‌های دیگر تابع خطای فوق را افزایش دهد. پس از تکمیل الگوریتم درختی برای نمونه‌های موجود در هر گره داخلی، یک مدل رگرسیون خطی چند متغیره برازش داده می‌شود. شکل ۱ نحوه تقسیم فضای مسئله به اجزای کوچک تر و کاربرد مدل‌های خطی را برای هر کدام از زیر دامنه‌ها با توجه به ساختار درختی نشان می‌دهد (Etemad-Shahidi and Taghipour, 2012).

در این تحقیق از سه مدل داده کاوی نوین استفاده شده است. الگوریتم M5 اولین بار توسط Quinlan (1992) ابداع شد و پس از آن توسط Wang and Witten (1997) توسعه و بهبود یافت. این مدل شاخه‌های خود را به صورت دوتایی و تنها بر اساس یک متغیر ایجاد می‌کند، بدین گونه که بر اساس شرطی که در هر گره تعریف می‌شود، اطلاعات در آن گره به دو قسمت تقسیم می‌شود (جباری و صمدی، ۱۳۹۲). در روش M5 فضای مسئله به زیر دامنه‌هایی تقسیم شده و برای هر زیر دامنه یک مدل رگرسیون خطی چند متغیره برازش داده می‌شود. این الگوریتم جداسازی‌های ممکن را در فضای چند متغیره انجام داده و به طور خودکار مدل‌هایی را برای هر کدام از دامنه‌ها ایجاد می‌کند (Quinlan, 1992). در این الگوریتم از پارامتر انحراف معیار مقادیر متغیر هدف به عنوان معیار اندازه‌گیری خطا در هر گره و ایجاد شاخه در آن گره استفاده می‌شود. بدین صورت که صفتی که موجب کاهش بیشتر انحراف معیار در هر گره می‌شود به عنوان صفت مورد نظر جهت ایجاد شاخه انتخاب می‌شود. کاهش انحراف استاندارد<sup>۱</sup> (SDR) که به عنوان تابع خطا در الگوریتم M5 به کار می‌رود به صورت زیر است:

$$SDR = sd(T) - \sum \left| \frac{T_i}{T} \right| \times sd(T) \quad (5)$$

در رابطه بالا T شامل نمونه‌هایی است که به گره مورد نظر رسیده‌اند و  $T_i$  شامل نمونه‌های است که از تقسیم گره مورد نظر بر اساس صفت انتخاب شده به دست آمده‌اند. Sd نیز معرف انحراف معیار است.

<sup>1</sup> Standard Deviation Reduction



شکل (۱): تقسیم فضای مسئله و ارائه مدل‌های خطی به ازای هر کدام از فضاها توسط الگوریتم M5

رشد خود رسیده سپس اصلاح می‌شوند. هدف مدل CART ایجاد یک درخت اصلاح شده نیست بلکه تولید یک سری درختان اصلاح شده تودرتو است که در آن همه درختان بهینه هستند. در درخت CART مجموعه‌ای از سؤالات به شکل  $x \leq d$  است که در آن  $x$  یک متغیر مستقل و  $d$  یک مقدار ثابت بوده و جواب هر سؤال آری یا نه است. در این روش در گره انتهایی مقدار عددی به‌عنوان نتیجه مورد نظر جهت آن گره ارائه می‌شود (Breiman et al., 1984). معیارهای متفاوتی جهت ایجاد شاخه و تولید درخت‌ها در روش CART وجود دارد که در تحقیق پیش رو از معیار انحراف حداقل مربعات<sup>۲</sup> (LSD) استفاده شده است. این معیار به‌صورت زیر تعریف می‌شود:

$$SS(t) = \sum_{i=1}^{N_t} (y_i(t) - \bar{y}(t))^2 \quad (6)$$

در رابطه ۶،  $N_t$  تعداد رکوردها،  $y_i(t)$  مقدار متغیر هدف و  $\bar{y}(t)$  میانگین مقادیر متغیر هدف در گره مورد نظر است. در الگوریتم CART یک متغیر ورودی زمانی به‌عنوان بهترین صفت برای ایجاد شاخه در گره  $t$  مورد استفاده قرار می‌گیرد که تابع زیر بیشینه شود.

$$Q(X, t) = SS(t) - [SS(t_R) + SS(t_L)] \quad (7)$$

برنامه‌ریزی ژنتیک<sup>۱</sup> (GP) از زیرشاخه‌های الگوریتم ژنتیک و از روش‌های الگوریتم گردش است که نخستین بار توسط Koza (1992) معرفی شد. برنامه‌ریزی ژنتیک از ساختار درختی برای بهینه‌سازی مسئله استفاده می‌کند، درحالی‌که الگوریتم ژنتیک بر روی رشته‌های بیتی کار می‌کند. در این روش هیچ ساختار و ارتباطی میان متغیرهای ورودی و خروجی وجود ندارد و ساختار مدل و ضرایب بهینه نیز طی فرآیند بهینه‌سازی به دست می‌آیند. ساختار درختی از مجموعه توابع (عملگرهای ریاضی مورد استفاده در فرمول‌ها) و ترمینال‌ها (متغیرهای مسئله و اعداد ثابت) ایجاد می‌شود (Koza, 1992). سه عمل ژنتیکی تلاقی، جهش و تولیدمثل از جمله مهم‌ترین عمل‌های ژنتیکی در برنامه‌ریزی ژنتیک می‌باشند. عمل‌های دیگر از قبیل اصلاح ساختار با احتمال کمتری بکار گرفته می‌شوند.

الگوریتم CART که اولین بار توسط Breiman et al. (1984) معرفی شد، به‌صورت یک درخت مرتبه‌ای دودویی است که فضای مسئله را به قسمت‌های جزء تقسیم می‌کند. در این روش داده‌ها به‌صورت خام و دست‌نخورده برای الگوریتم معرفی شده و هیچ‌گونه فیلتری بروی داده‌ها صورت نمی‌گیرد. در این مدل درختان بدون وجود عامل متوقف‌کننده به حداکثر

<sup>2</sup> Least-Squared Deviation

<sup>1</sup> Genetic Programming

سال هشتم • شماره بیست و نهم • پاییز ۱۳۹۶

در روابط فوق  $K_m$  ضریب انتشار طولی مشاهداتی و  $K_c$  ضریب انتشار طولی محاسباتی است.

### نتایج و بحث

پس از بررسی‌های آماری صورت گرفته بر روی روابط تجربی، سه رابطه (Seo and Cheong (1998)، Rajeev و Kashefipour and Falconer (2002) and Dutta (2009) به‌عنوان بهترین روابط در پیش‌بینی ضریب پراکندگی طولی در این تحقیق انتخاب و برای مقایسه با سه مدل داده‌کاوی به کار گرفته شدند. شکل ۲ عملکرد سه رابطه تجربی نام‌برده را برای تخمین  $K_x$  نشان می‌دهد.

در معادله ۷،  $SS(t)$  مجموع مربعات خطا و  $SS(t_R)$  مجموع مربعات خطا به ترتیب در شاخه سمت راست و چپ گره  $t$  است.

مدل M5 تنها قادر به شبیه‌سازی مدل‌های خطی بوده که این امر با اساس معادلات متداول ضریب انتشار طولی در تناقض است. بر همین اساس کلیه داده‌ها به صورت لگاریتم طبیعی به مدل M5 معرفی شدند و پس از ساخت مدل خطی از حالت لگاریتمی به توانی تبدیل شدند. برای مدل‌های GP و CART هم در ابتدا از لگاریتم طبیعی اعداد استفاده شد و پس از ساخت مدل، خروجی‌ها از حالت لگاریتمی به حالت توانی تبدیل شدند تا شرایط کاملاً یکسانی برای هر سه مدل داده‌کاوی فراهم شده باشد. در تحقیق حاضر از مجموع اطلاعات در دسترس ۸۰٪ (۱۲۰ سری از داده‌ها) جهت آموزش مدل‌ها و ۲۰٪ باقیمانده (۲۹ سری از داده‌ها) جهت صحت‌سنجی مورد استفاده قرار گرفتند. برای مدل‌سازی M5، GP و CART به ترتیب از برنامه‌های WEKA 3.7، STATISTICA و GeneXproTools 5.0 استفاده شد. جهت مقایسه بهتر کارایی مدل‌های ارائه شده از آنالیزهای آماری خطای متوسط مطلق<sup>۱</sup> (MAE)، جذر میانگین مربعات خطا<sup>۲</sup> (RMSE) و ضریب تبیین<sup>۳</sup> ( $R^2$ ) استفاده گردید. علاوه بر معیارهای ذکر شده در این تحقیق از نسبت اختلاف (DR) نیز استفاده شده است. این نسبت توسط White et al. (1973) پیشنهاد شده است و از معیارهای قدرتمند آماری به حساب می‌آید.

$$MAE = \frac{1}{N} \sum |D_{Ri}| \quad (8)$$

$$RMSE = \sqrt{\frac{1}{N} \sum (D_{Ri})^2} \quad (9)$$

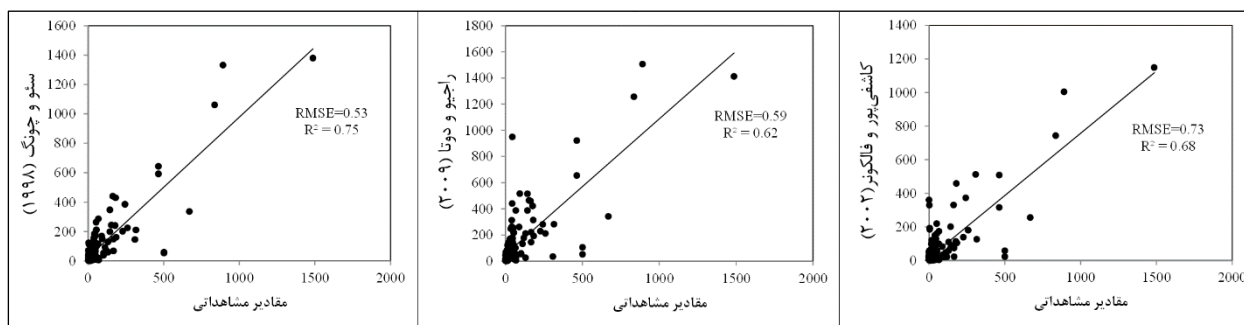
$$R^2 = 1 - \frac{\sum (K_m - K_c)^2}{\sum (K_m - \bar{K}_m)^2} \quad (10)$$

$$D_R = \log \frac{K_c}{K_m} \quad (11)$$

<sup>1</sup> Mean-Absolute Error

<sup>2</sup> Root Mean Square Error

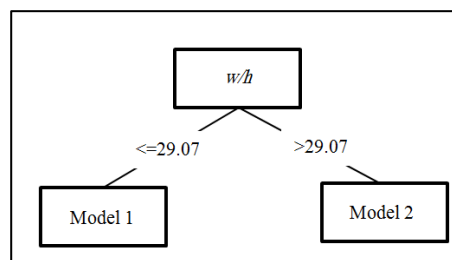
<sup>3</sup> Determination coefficient



شکل (۲): عملکرد روابط تجربی انتخاب شده در تخمین ضریب پخشیدگی طولی

با توجه به شکل ۳، پارامتر مورد استفاده جهت تشکیل مدل درختی، نسبت عرض آبراهه به عمق جریان  $(w/h)$  می‌باشد. از لحاظ فیزیکی باریک یا عریض بودن رودخانه تاثیر زیادی بر روی میزان پخشیدگی طولی داشته و  $(w/h) \leq 29.07$  می‌تواند نماینده رودخانه‌های باریک و  $(w/h) > 29.07$  نشان‌دهنده رودخانه‌های عریض باشد. با توجه به معادلات ۱۲ و ۱۳، در آبراهه‌هایی که نسبت عرض به عمق آن‌ها کمتر از  $29/07$  باشد، ترم  $(w/h)$  تأثیر بیشتری نسبت به ترم  $(u/u_*)$  در پیش‌بینی ضریب پراکندگی طولی خواهد داشت. در صورتی که در آبراهه‌های عریض که نسبت عرض به عمق آن‌ها بیش از  $29/07$  است، شاهد افزایش تأثیر ترم  $(u/u_*)$  نسبت به  $(w/h)$ ، در تخمین ضریب پراکندگی خواهیم بود. در رودخانه‌های باریک با  $(w/h)$  کوچک تأثیر تنش برشی ناچیز است ولی با افزایش  $(w/h)$  تأثیر تنش برشی بر روی ضریب پخشیدگی افزایش می‌یابد (Papadimitrakis and Orphanos, 2004). علاوه بر این Rutherford (1994) نیز در مطالعه‌ای به این نتیجه رسیده است که در رودخانه‌های عریض نسبت به رودخانه‌های باریک، پارامتر سرعت تأثیر بیشتری بر روی ضریب پخشیدگی دارد. شکل‌های ۴ و ۵ نمایش‌دهنده نحوه عملکرد مدل M5 در مرحله آموزش و صحت سنجی برای تخمین ضریب پراکندگی طولی می‌باشند.

پس از معرفی داده‌های مشاهده‌ای به صورت لگاریتم طبیعی به M5، مدل مورد آموزش قرار گرفت. دو معادله ارائه شده توسط این مدل، روابط خطی بوده که توسط الگوریتم درختی برای هر کدام از حالت‌ها توسعه یافته است. نمودار درختی مدل پیشنهادی M5 در شکل ۳ ارائه شده است.

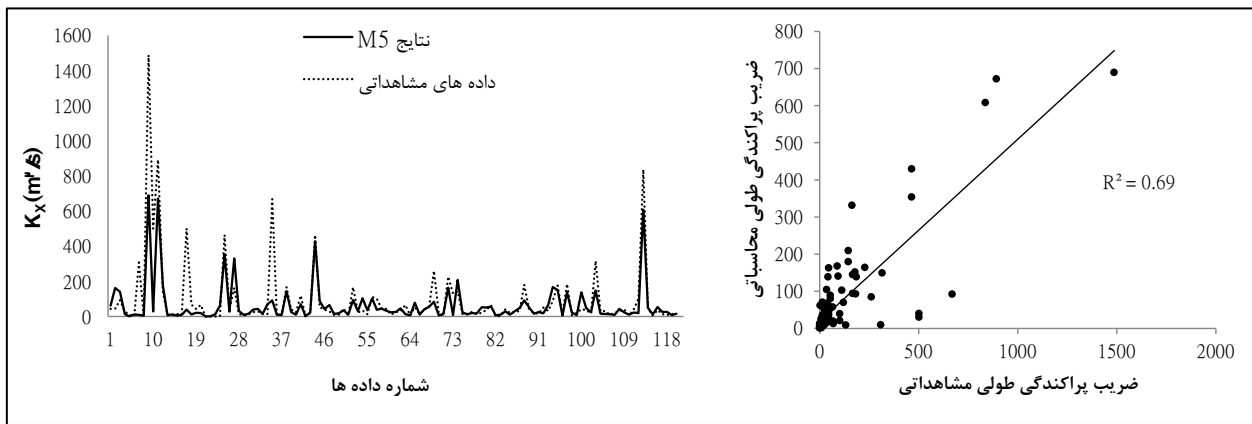


شکل (۳): ساختار درختی M5 جهت برآورد ضریب پراکندگی طولی در رودخانه‌ها

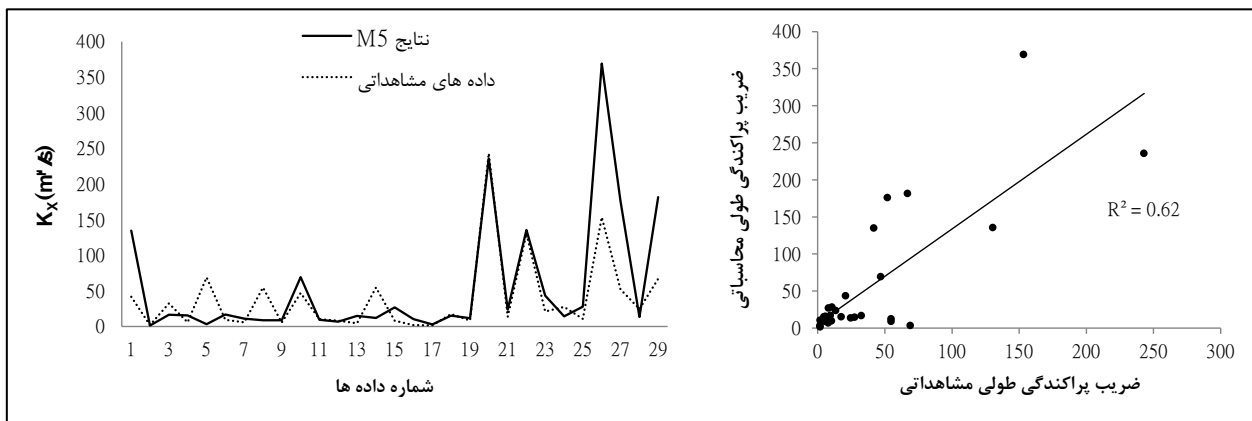
روابط ارائه شده توسط مدل M5 پس از تبدیل از حالت لگاریتمی، به صورت معادلات ۱۲ و ۱۳ می‌باشند.

$$Model 1: IF \left( \frac{w}{h} \right) \leq 29.07 \Rightarrow \left( \frac{K_x}{hu_*} \right) = 12.06 \left( \frac{w}{h} \right)^{0.8} \left( \frac{u}{u_*} \right)^{0.15} \quad (12)$$

$$Model 2: IF \left( \frac{w}{h} \right) > 29.07 \Rightarrow \left( \frac{K_x}{hu_*} \right) = 13.73 \left( \frac{w}{h} \right)^{0.62} \left( \frac{u}{u_*} \right)^{0.91} \quad (13)$$



شکل (۴): عملکرد مدل M5 در مرحله آموزش برای تخمین ضریب پراکندگی طولی



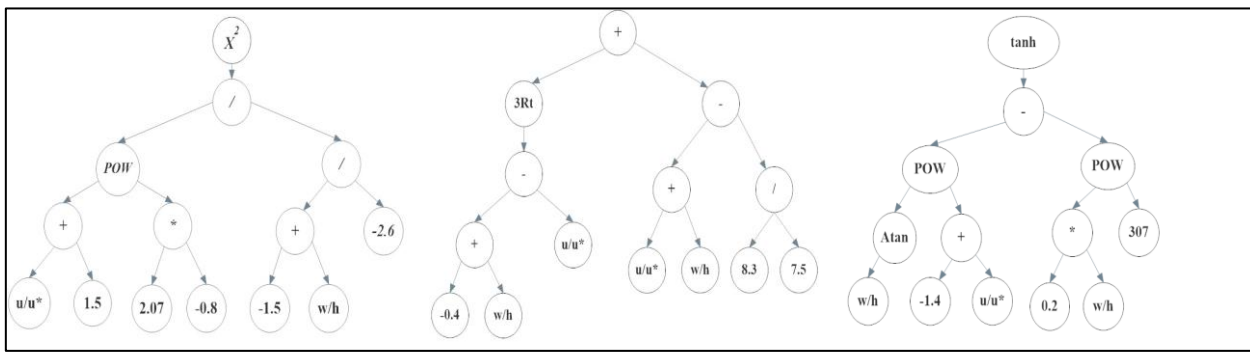
شکل (۵): عملکرد مدل M5 در مرحله صحت سنجی برای تخمین ضریب پراکندگی طولی

گردید. معادله ۱۳ ضمن پیچیدگی فراوان، به عنوان بهترین و دقیق ترین مدل ارائه شده توسط مدل GP انتخاب شد. شکل ۶ نشان دهنده ساختار درختی مدل GP جهت برآورد ضریب پراکندگی طولی بوده که حاصل جمع این سه درخت، مدل نهایی را ارائه می دهد. شکل های ۷ و ۸ نمایش دهنده نحوه عملکرد مدل GP در مرحله آموزش و صحت سنجی برای تخمین ضریب پراکندگی طولی است.

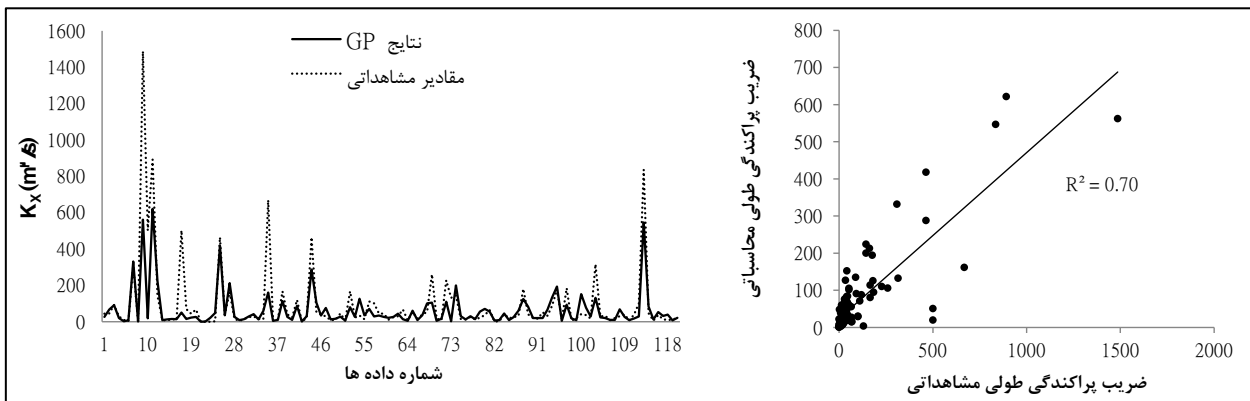
جهت مقایسه دقیق تر کارایی الگوریتم های درختی، داده های ورودی در مراحل آموزش و صحت سنجی برای تمامی مدل ها یکسان معرفی شدند. علاوه بر آن داده های ورودی در ابتدا برای تمامی مدل ها به صورت لگاریتم طبیعی بوده و پس از اجرای مدل ها، نتایج به مقادیر واقعی تبدیل شدند. در مدل GP ابتدا از چهار عملگر اصلی جمع، تفریق، ضرب و تقسیم استفاده شد، و با افزودن عملگرهای tan<sup>-1</sup> و tanh، power، sqrt نتایج بهتری حاصل

$$\left( \frac{K_x}{hu_*} \right) = \exp \left( \left( \frac{\left( \frac{u}{u_*} + 1.5 \right)^{-1.6}}{\left( \frac{-1.5 + \frac{w}{h}}{-2.6} \right)} \right)^2 + \left( \sqrt[3]{\left( -0.4 + \frac{w}{h} \right) - \frac{u}{u_*} + \frac{w}{h} + \frac{u}{u_*} + 20.75} \right) + \left( \tanh \left( \tan^{-1} \left( \frac{w}{h} \right) \right) \right)^{-1.4 + \frac{u}{u_*}} - \left( \left( 0.2 * \frac{w}{h} \right)^{307} \right) \right) \quad (13)$$

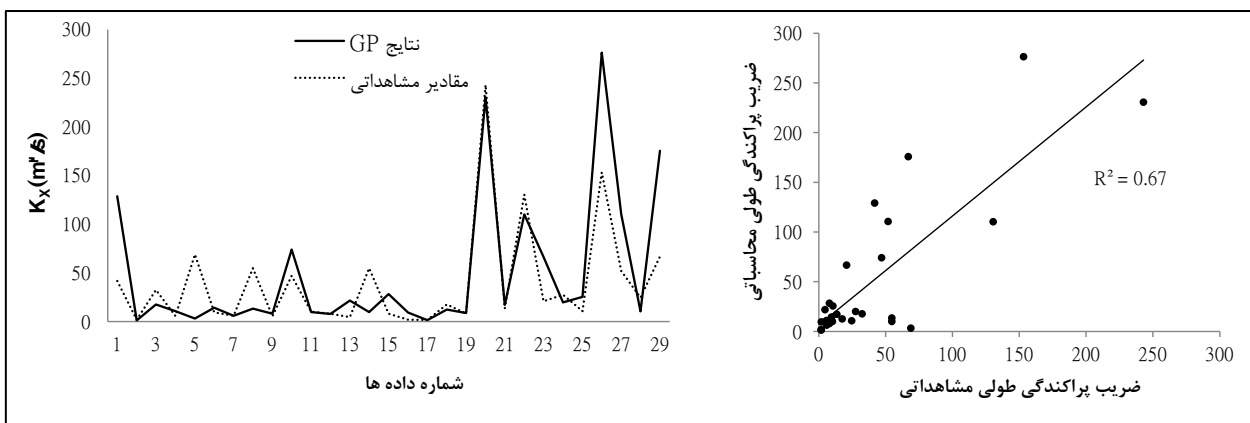




شکل (۶): ساختار درختی GP جهت برآورد ضریب پراکندگی طولی



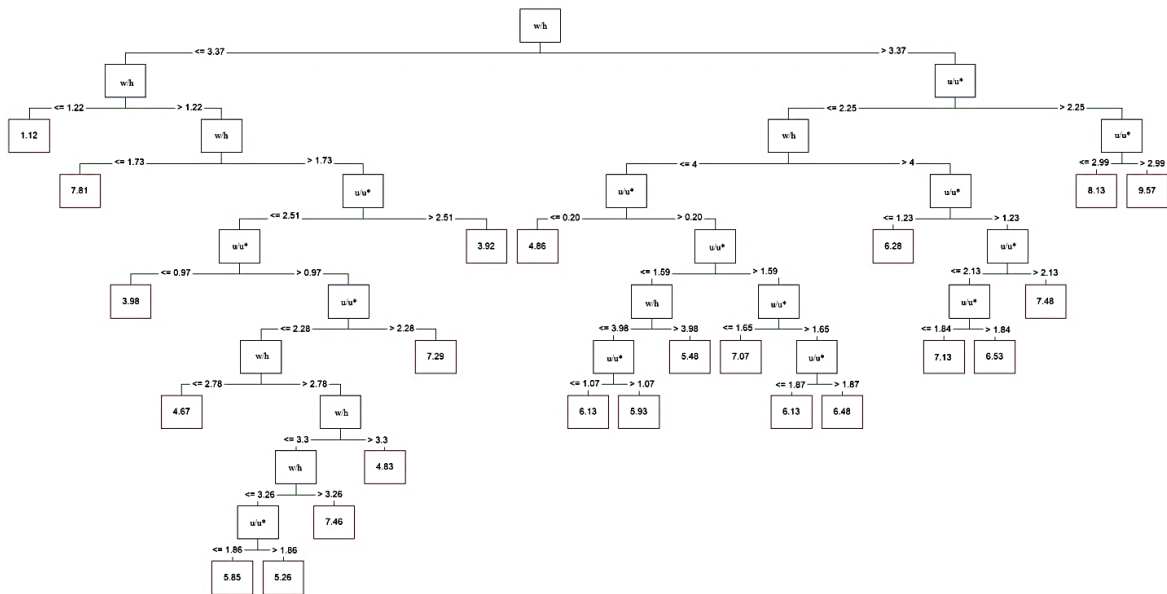
شکل (۷): عملکرد مدل GP در مرحله آموزش برای تخمین ضریب پراکندگی طولی



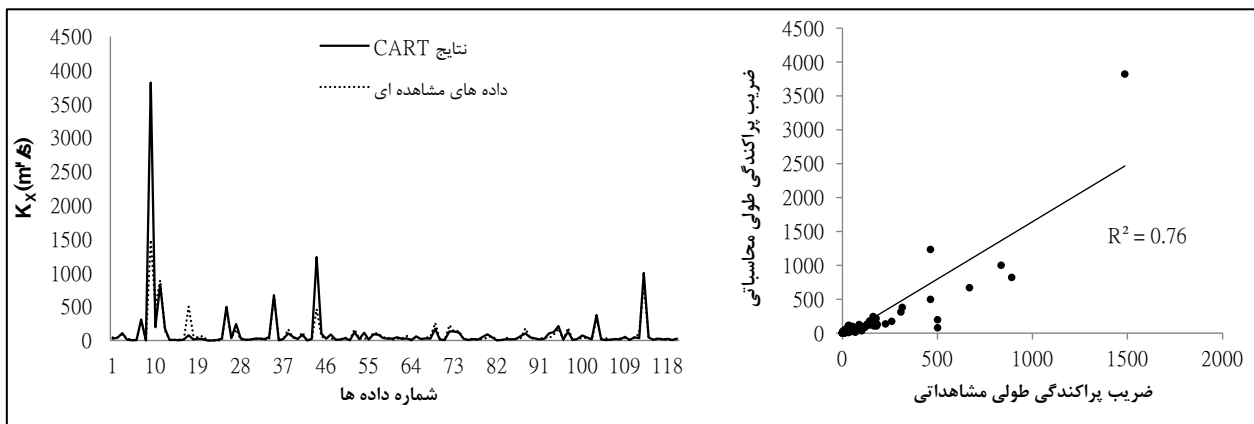
شکل (۸): عملکرد مدل GP در مرحله صحت سنجی برای تخمین ضریب پراکندگی طولی

CART از گره های متعددی تشکیل شده و در انتهای هر گره یک کمیت عددی به عنوان خروجی ارائه می دهد (Avg). این کمیت بیانگر میانگین ضریب پراکندگی طولی مربوط به اطلاعات آن گره است. تعداد زیاد گره ها در این الگوریتم باعث مشکل شدن کاربرد آن در سایر مدل ها می گردد. همچنین اشکال

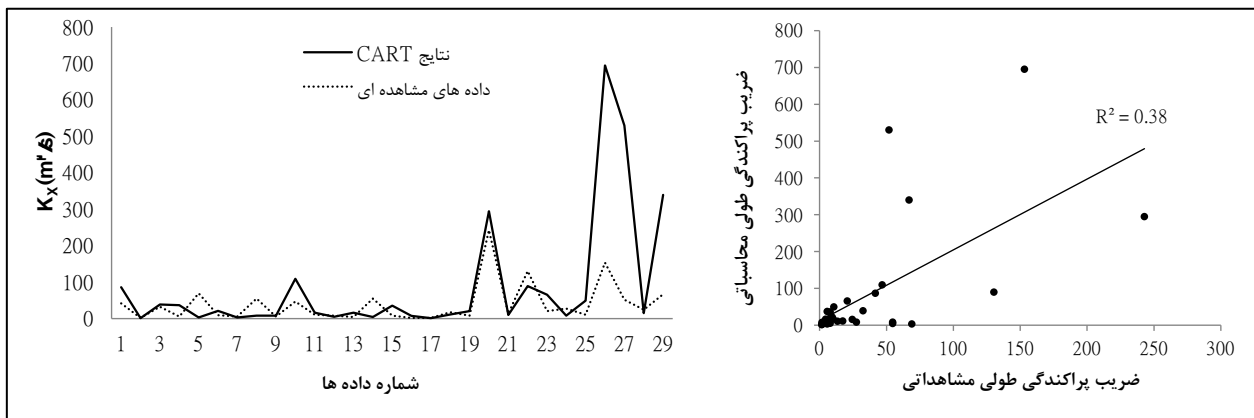
خروجی مدل درختی CART بر اساس داده های ورودی در شکل ۹ نشان داده شده است. همچنان که در این شکل مشاهده می شود، در اولین گره، ترم  $(w/h)$  به عنوان بهترین صفت جهت ایجاد شاخه انتخاب گردید که این امر نشان دهنده تأثیر شکل هندسی آبراهه بر ضریب پراکندگی طولی است. همان گونه که در شکل ۹ مشاهده می شود، الگوریتم



شکل (۹): ساختار درختی CART جهت برآورد ضریب پراکندگی طولی



شکل (۱۰): عملکرد مدل CART در مرحله آموزش برای تخمین ضریب پراکندگی طولی

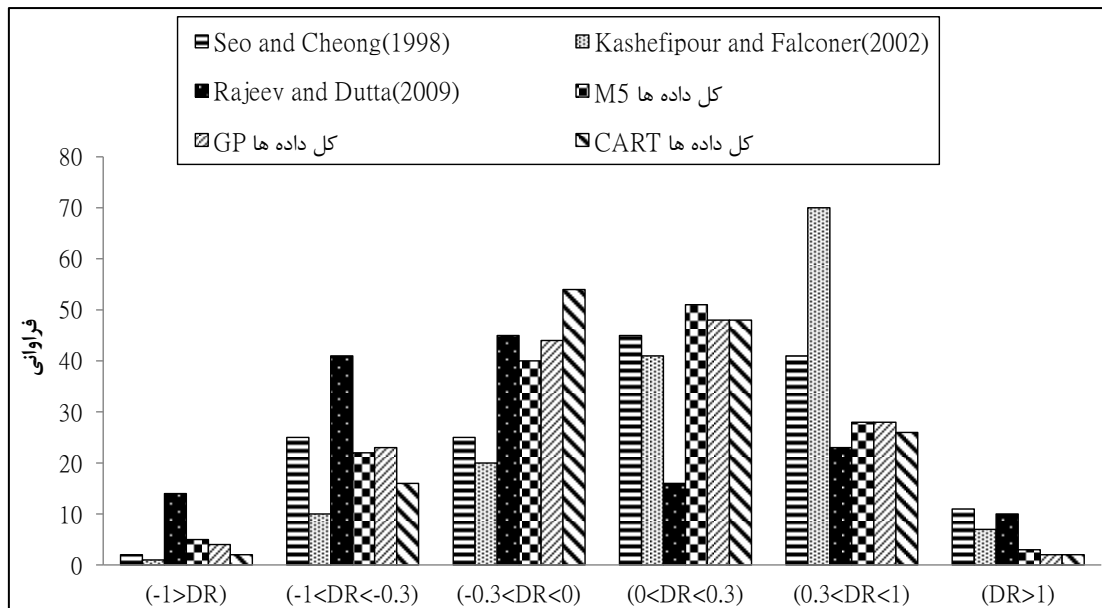


شکل (۱۱): عملکرد مدل CART در مرحله صحت سنجی برای تخمین ضریب پراکندگی طولی

جهت بررسی دقت مدل‌ها از شاخص‌های آماری متنوعی استفاده گردید که یکی از مهم‌ترین شاخص‌ها، DR است. چنانچه DR برابر صفر باشد نشان‌دهنده تطابق کامل نتایج مدل‌ها و مقادیر واقعی است. در غیر این صورت چنانچه  $(DR > 0)$  باشد مدل ضریب پراکندگی طولی را بیشتر از مقدار واقعی<sup>۱</sup> برآورد کرده است در غیر این صورت  $(DR < 0)$ ، ضریب پراکندگی طولی کمتر از مقدار واقعی محاسبه شده است. در این تحقیق درصد مقادیر با DR بین  $0/3 -$  و  $0/3$  به‌عنوان دقت<sup>۲</sup> هر مدل در نظر گرفته شده است که مطابق دقت تعریف شده توسط Seo and Cheong (1998) است. هیستوگرام فراوانی DR برای سه روش تجربی و مدل‌های درختی در شکل ۱۲ نشان داده شده است. علاوه بر این نتایج تحلیل‌های آماری صورت گرفته در جدول ۳ ارائه شده است.

---

<sup>1</sup> Overestimation<sup>2</sup> Accuracy



شکل (۱۲): مقایسه نسبت اختلاف (DR) مدل‌ها

جدول (۳): بررسی کارایی مدل‌های هوشمند و تجربی در تخمین ضریب پخشیدگی طولی

R <sup>2</sup>	RMSE	MAE	Accuracy %	DR<-1 %	-1<DR<-0.3 %	-0.3<DR<0 %	0<DR<0.3 %	0.3<DR<1 %	DR >1 %	مدل
۰/۷۵	۰/۵۳	۰/۴۳	۴۶/۹	۱/۳	۱۶/۷	۱۶/۷	۳۰/۲	۲۷/۵	۷/۳	Seo and Cheong(1998)
۰/۶۸	۰/۷۳	۰/۵۳	۴۰/۹	۹/۳	۲۷/۵	۳۰/۲	۱۰/۷	۱۵/۴	۶/۷	Kashefipour and Falconer(2002)
۰/۶۲	۰/۵۹	۰/۴۷	۴۰/۹	۰/۶	۶/۷	۱۳/۴	۲۷/۵	۴۶/۹	۴/۶	Rajeev and Dutta(2009)
۰/۶۵	۰/۴۴	۰/۳۲	۶۱	۳/۳	۱۴/۷	۲۶/۸	۳۴/۲	۱۸/۷	۲	کل داده‌ها
۰/۶۲	۰/۴۴	۰/۳۵	۵۵	۳/۴	۶/۸	۲۷/۵	۲۷/۵	۳۴/۴	۰	M5 صحت سنجی
۰/۶۸	۰/۴۱	۰/۲۹	۶۱/۷	۲/۶	۱۵/۴	۲۹/۵	۳۲/۲	۱۸/۷	۱/۳	کل داده‌ها
۰/۶۷	۰/۴۳	۰/۳۱	۵۸/۵	۳/۴	۱۰/۳	۳۱	۲۷/۵	۲۷/۵	۰	GP صحت سنجی
۰/۷۴	۰/۳۶	۰/۲۵	۶۸/۴	۱/۳	۱۰/۷	۳۶/۲	۳۲/۲	۱۷/۴	۱/۳	کل داده‌ها
۰/۳۸	۰/۵۶	۰/۴۶	۴۱/۲	۶/۸	۶/۸	۲۷/۵	۱۳/۷	۴۱/۳	۳/۴	CART صحت سنجی

### نتیجه‌گیری

هدف از این تحقیق برآورد ضریب پخشیدگی طولی در آبراه‌ها با استفاده از مدل‌های داده‌کاوی M5، GP و CART می‌باشد. طبق نتایج به‌دست آمده، روابط تجربی از دقت پایین‌تری در مقایسه با مدل‌های داده‌کاوی برخوردارند. در میان مدل‌های داده‌کاوی مورد استفاده، مدل CART در مرحله آموزش و صحت‌سنجی به ترتیب دارای ضریب تبیین ۰/۷۶ و ۰/۳۸ بوده که خود بیانگر قدرت بالای این مدل در مرحله آموزش و ضعف آن در مرحله صحت‌سنجی است. دقت مدل CART برای کل داده‌ها و داده‌های صحت‌سنجی به ترتیب ۶۸/۴ و ۴۱/۲ است. بررسی نتایج سایر مدل‌ها نشان می‌دهد که مدل پیشنهادی CART در مرحله صحت‌سنجی ضعیف ظاهر شده است. در مقابل نتایج آماری دو مدل GP و M5 نزدیک به هم بوده اما رابطه پیشنهادی توسط مدل GP به مراتب پیچیده‌تر از دو رابطه پیشنهادی M5 است. بر اساس آنالیزهای صورت گرفته مدل پیشنهادی M5 دارای دقت ۶۱ و ۵۵ درصد به ترتیب برای کل داده‌ها و داده‌های صحت‌سنجی است که این میزان دقت به طرز چشم‌گیری بیشتر از دقت ارائه شده توسط روابط تجربی است. بر این اساس می‌توان از مدل M5 با توجه به دقت بالا و سادگی معادلات ارائه شده و قابل درک بودن آن جهت برآورد ضریب پراکندگی طولی استفاده کرد.

با توجه به شکل ۱۲ و جدول ۳، دقت مدل‌های درختی به کار گرفته شده همگی بیشتر از ۶۰ درصد بوده است. این در حالی است که رابطه تجربی Seo and Cheong (1998) با توجه به دارا بودن شاخص‌های آماری مطلوب‌تر نسبت به سایر روابط تجربی، تنها دارای دقتی معادل ۰/۴۶ است که این امر بیانگر دقت بیشتر مدل‌های درختی نسبت به روابط تجربی در تخمین ضریب پراکندگی طولی است. در میان مدل‌های درختی، مدل CART برای کل داده‌ها با دارا بودن ضریب تبیین ۰/۷۴ و جذر میانگین مربعات خطا ۰/۳۶، از دقت بیشتری نسبت به دیگر مدل‌ها برخوردار است. این در حالی است که دقت این مدل در مرحله صحت‌سنجی بسیار پایین بوده که نشان‌دهنده ضعف این مدل در تعمیم نتایج است. نزدیک بودن نتایج تحلیل آماری دو مدل پیشنهادی GP و M5 در مراحل آموزش و صحت‌سنجی بیانگر آموزش صحیح این دو الگوریتم است. دقت دو مدل GP و M5 بسیار نزدیک به هم بوده، اما پیچیدگی فراوان مدل GP و سادگی و قابل درک بودن مدل M5 نشان‌دهنده برتری این مدل در تخمین ضریب پخشیدگی طولی است. در مدل M5، جذر میانگین مربعات خطا که در تحقیقات مربوط به ضریب پراکندگی طولی کاربرد فراوانی دارد، نسبت به رابطه Seo and Cheong (1998) از ۰/۵۳ به ۰/۴۴ کاهش یافته که این امر نشان‌دهنده کارایی این الگوریتم در تخمین ضریب پراکندگی طولی است.

### منابع

ظهیری، ج. ۱۳۹۴. کاربرد مدل‌های ناپارامتریک CART و M5 در محاسبه عمق آبستگي اطراف پایه‌های پل. فصلنامه علمی پژوهشی مهندسی آبیاری و آب، شماره ۲۰، ص ۳۵-۵۰.

Breiman, L., Friedman, J., Olshen, R. and Stone, C. 1984. Classification and Regression Trees. Chapman & Hall/CRC Press, Boca Raton, FL.

Chatila, G. J. 1997. Modeling of pollutant transfer in compound open channels. PhD Dissertation, University of Ottawa, Ontario, Canada.

سال هشتم • شماره بیست و نهم • پاییز ۱۳۹۶

Deng, Z.Q., Bengtsson, L., Singh, V. P., et al. 2002. Longitudinal dispersion coefficient in single-channel streams. *Journal of Hydraulic Engineering*. 128 (10): 901-916.

Etemad-Shahidi, A and M. Taghipour. 2012. Predicting longitudinal dispersion coefficient in natural streams using M5' model tree. *Journal of Hydraulic Engineering*, 138(6): 542-554.

Etemad-Shahidi, A and N. Ghaemi. 2011. Model tree approach for prediction of pile groups scour due to waves. *Ocean Engineering*, 38: 1522-1527.

Haghiabi, A. H. 2016. Prediction of longitudinal dispersion coefficient using multivariate adaptive regression splines. *Journal of Earth System Science*, 125: 985-995.

Kashefipour, M.S., Falconer, R.A. 2002. Longitudinal dispersion coefficients in natural channels. *Water Res.* 36 (6): 1596-1608.

Koza, J.R. 1992 . *Genetic Programming: on the programming of computers by means of natural selection.* Cambridge, MA: MIT Press.

Li, Z.H., Huang, J., Li, J. 1998. Preliminary study on longitudinal dispersion coefficient for the gorges reservoir. In: *Proceedings of the Seventh International Symposium Environmental Hydraulics*, 16-18 December, Hong Kong, China.

Noori, R., Karbassi, A., Farokhnia, A and Dehghani, M. 2009. Predicting the longitudinal dispersion coefficient using support vector machine and adaptive Neuro-Fuzzy inference system techniques; *Environmental Engineering Science*. 26(10):1503-1510.

Papadimitrakis, I., and Orphanos, I. 2004. Longitudinal dispersion characteristics of rivers and natural streams in Greece. *Water, Air, & Soil Pollution*, 4(4-5): 289-305.

Quinlan, J. R. 1992. *Learning with continuous classes.* Proc., 5<sup>th</sup> Australian Joint Conf. on Artificial Intelligence, World Scientific, Singapore, 343-348.

Rutherford, J.C. 1994. *River Mixing.* John Wiley, Chichester, U. K.

Sattar, A. M. A and Gharabaghi B .2015. Gene expression models for prediction of longitudinal dispersion coefficient in streams; *Journal of Hydrology*. 524: 587-596.

Seo, I and Baek, K .2004 .Estimation of the longitudinal dispersion coefficient using the velocity profile in natural streams; *Journal of Hydraulic Engineering*, 130(3): 227-236.

Seo, I and Cheong, T. 1998 .Predicting longitudinal dispersion coefficient in natural streams; *Journal of Hydraulic Engineering*, 124(1): 25-32.

Tayfur, G., Singh, V.P. 2005. Predicting longitudinal dispersion coefficient in natural streams by artificial neural network. *Journal of Hydraulic Engineering*, 131 (11): 991-1000.

Wang, Y., and Witten, I. H. 1997. Induction of model trees for predicting continuous classes. Proc. of the Poster Papers of the European Conf. on Machine Learning, Univ. of Economics, Faculty of Informatics and Statistics, Prague, Czech Republic.

## Investigation of Tree Models Performance for Estimation of Longitudinal Dispersion Coefficient in Straight River

Hosein Nezaratian<sup>1</sup>, Javad Zahiri<sup>2</sup>, Seyed Mahmood Kashefipour<sup>3</sup>

### Abstract

Modeling pollution transmission in rivers is an important subject in environmental studies. Longitudinal dispersion coefficient is one of the key factors in the modelling of lateral dispersion of pollutants. Several researchers have attempted to estimate this coefficient using empirical and semi-empirical methods. However, robust models that can accurately estimate longitudinal dispersion coefficient in river streams are still required. In this study, data driven models were developed using the hydraulic and geometric parameters of rivers. The classification and regression tree (CART), M5 and genetic programming (GP) were used for this purpose. The models performances were then compared quantitatively with those of existing ones using accuracy parameters such as root mean square error (RMSE), mean absolute error (MAE) and discrepancy ratio (DR). The results illustrated that data driven models outperform the existing formulae in term of accuracy. CART model outperform other models in training step, but its performance decrease for testing data. M5 and GP models have RMSE of 0.41 and 0.44 and accuracy of 61% and 62%, respectively. According to small difference between M5 and GP performances, and simple structure of M5 algorithm, this model can be used for estimating longitudinal dispersion coefficient in streams.

**Keywords:** Pollution Transmission, Longitudinal Dispersion Coefficient, CART, GP, M5

---

<sup>1</sup> M. Sc., Water structures, Khuzestan Ramin Agriculture and Natural Resources University, Hosein.nezaratian69@gmail.com

<sup>2</sup> Assistant Professor, Department of Water Engineering, Khuzestan Ramin Agriculture and Natural Resources University, Zahiri\_Javad@yahoo.com. (Corresponding author)

<sup>3</sup> Professor, Department of Water Engineering, Shahid Chamran University of Ahvaz, Kashefipour@excite.com