



مقایسه روش‌های K- نزدیک‌ترین همسایگی و شبکه عصبی مصنوعی در برآورد ظرفیت تبادل کاتیونی خاک

*علی‌اصغر ذوالفقاری^۱، محمدتقی تیرگر سلطانی^۲، تورج افشاری بدرلو^۳ و فریدون سرمدیان^۴
^۱ دانشجوی دکتری گروه مهندسی علوم خاک، دانشگاه تهران، کارشناس ارشد گروه مهندسی علوم خاک، دانشگاه تهران،
^۲ دانشجوی کارشناسی ارشد گروه مهندسی علوم خاک، دانشگاه صنعتی شاهرود، ^۳ دانشیار گروه مهندسی علوم خاک، دانشگاه تهران
تاریخ دریافت: ۹۱/۳/۹؛ تاریخ پذیرش: ۹۱/۹/۲۵

چکیده

اندازه‌گیری ظرفیت تبادل کاتیونی خاک در سطوح وسیع، معمولاً بسیار پرهزینه و وقت‌گیر است. تخمین این کمیت به‌وسیله ویژگی‌های زودیافت خاک، از طریق توسعه توابع غیرپارامتریک می‌تواند رویکرد مناسبی باشد. در این پژوهش روش غیرپارامتریکی با عنوان K- نزدیک‌ترین همسایگی در تخمین CEC خاک استفاده شد و نتایج آن با یکی از پرکاربردترین روش‌های مرسوم مبتنی بر مدل‌های شبکه عصبی مصنوعی (ANN) مورد مقایسه قرار گرفت. ۶۸۳ نمونه خاک از مناطق مرکزی ایران انتخاب شدند که ۱۲۰ عدد از آن‌ها به‌عنوان داده‌های مورد آزمون (هدف) و ۵۶۳ عدد به‌عنوان بانک داده مرجع (آموزش) قرار گرفتند. مقادیر پارامترهای رس، سیلت، شن و کربن آلی خاک به‌عنوان متغیر مستقل ورودی (زودیافت) و CEC به‌عنوان متغیر وابسته خروجی بودند. نتایج نشان داد که بیش‌ترین خطای برآورد (MaxE) در روش K-NN برابر $4/81 \text{ cmol}^+/\text{kg}$ و این مقدار در روش ANN برابر $5/26 \text{ cmol}^+/\text{kg}$ بود. ریشه میانگین مربعات خطا در روش K-NN، $1/51$ و در روش ANN، $1/53$ بود، که نشان می‌دهد هر دو روش قادرند با دقت بالا و یکسانی CEC خاک‌های هدف را پیش‌بینی نمایند. مقادیر مثبت آماره میانگین خطا (ME) برای این دو روش نیز نشان داد که هر دوی آن‌ها متمایل به برآورد کم‌تر مقدار CEC می‌باشند. همچنین نتایج بررسی کارایی مدل‌ها نشان داد که هر دو روش از کارایی بالایی ($EF=0/88$) در برآورد ظرفیت تبادل کاتیونی خاک برخوردار هستند.

واژه‌های کلیدی: مدل غیرپارامتریک، ظرفیت تبادل کاتیونی، K- نزدیک‌ترین همسایگی، مدل شبکه

عصبی مصنوعی

* مسئول مکاتبه: azolfaghari@ut.ac.ir

مقدمه

ظرفیت تبادل کاتیونی (CEC)^۱ یکی از ویژگی‌های مهم خاک است که بسیاری از عملکردهای اساسی خاک در ارتباط با نوسانات اسیدیته (pH)، نگهداری آب، مواد غذایی و آلاینده‌ها را کنترل می‌نماید و به همین دلیل این پارامتر، یکی از شاخص‌های کلیدی در مدیریت حاصل‌خیزی خاک و بهره‌وری زمین محسوب می‌شود (کارو و همکاران، ۲۰۰۱). تعیین سطوح بحرانی آلاینده‌های شیمیایی در خاک و پالایش آن‌ها از خاک، بدون دسترسی به داده‌های CEC امکان‌پذیر نمی‌باشد. داده‌های CEC خاک‌های زراعی می‌تواند تفسیر روشن‌تر و کامل‌تری از فرایندهای تغذیه گیاه، کاربرد کودها و اصلاح شیمیایی خاک‌ها را در اختیار پژوهشگران قرار دهد. دقیق‌ترین داده‌های CEC خاک‌ها در آزمایشگاه و از طریق اندازه‌گیری مستقیم به دست می‌آید، ولی روش‌های اندازه‌گیری این پارامتر به نسبت پرهزینه و البته بسیار وقت‌گیر می‌باشند (سیبلد و همکاران، ۲۰۰۵؛ مک‌براتی و همکاران، ۲۰۰۲). به همین دلیل تاکنون بسیاری از پژوهشگران بررسی‌های خود را بر روش‌های غیرمستقیم برآورد CEC از روی ویژگی‌های زودریافت خاک متمرکز نموده‌اند (سیبلد و همکاران، ۲۰۰۵). رگرسیون‌گیری بین ویژگی‌های پایه‌ای خاک به عنوان داده‌های ورودی و مقادیر CEC به دست آمده در آزمایشگاه، منجر به ایجاد معادله‌های چندمتغیره‌ای می‌شود که می‌تواند تخمین قابل‌قبولی از CEC خاک مناطق همگن را در اختیار پژوهشگران قرار دهد. این معادله‌ها که CEC را تابعی خطی از ویژگی‌های فیزیکی و شیمیایی خاک به خصوص مقدار رس و مواد آلی خاک تخمین می‌زنند تحت عنوان کلی توابع انتقالی خاک (PTFs)^۲ نام گرفته‌اند (مک‌براتی و همکاران، ۲۰۰۲). کریمیان (۱۹۹۸) سهم رس و مواد آلی در CEC ۱۵۰ نمونه از خاک‌های آهکی مناطقی استان فارس را در غالب توابع انتقالی بیان نمود. کروگ و همکاران (۲۰۰۰) از توابع انتقالی برای پیش‌بینی ظرفیت تبادل کاتیونی ۱۶۴۳ نمونه از خاک‌های دانمارک استفاده کردند. آن‌ها با متوسط خطای پیش‌بینی $1/99 \text{ cmol}^+/\text{kg}$ نشان دادند که ۹۰ درصد تغییرات CEC در خاک‌های مورد آزمایش، به وسیله مقدار رس و مواد آلی قابل توجیه بود. بل و وان‌کولن (۱۹۹۵) در پژوهشی در مکزیک تابعی ارایه نمودند که بیش از ۹۶ درصد تغییرات CEC را با مقادیر رس، ماده آلی و pH توجیه می‌نمود. از عیب‌های PTF ها و به طور کلی توابع پارامتریک در تخمین ویژگی‌های دیریافت خاک تعیین مقدار عددی هر یک از پارامترها

1- Cation Exchange Capacity

2- Pedotransfer Function

است که از طریق بهترین برازش معادله‌های یاد شده بر داده‌ها به دست می‌آید. از آنجا که مقادیر این پارامترها برای هر سری داده‌ها تغییر می‌کند، عیب‌های اصلی کاربرد این توابع این است که این معادله‌ها برای هر سری از داده‌های جدید نیازمند برازش دوباره و تعیین مقادیر جدید برای پارامترها است. ضمن آنکه تحقق شرایط اعتبار رگرسیون‌گیری (مانند نرمال بودن توزیع خطاها حول میانگین و...) از محدودیت‌های کاربرد این توابع می‌باشد. علاوه بر این تخمین نتایج در سری‌های کوچک داده‌ای دارای اریب است (نمس و همکاران، ۲۰۰۶). اخیراً پژوهشگران با استفاده از روش‌های غیرپارامتریک مانند مدل‌های شبکه عصبی مصنوعی (ANN)^۱ تخمین‌های دقیق‌تری از مقدار CEC ارائه نموده‌اند (معماریان‌فرد و بیگی‌هرچگانی، ۲۰۰۹؛ امینی و همکاران، ۲۰۰۵). شبکه‌های عصبی مصنوعی یکی از روش‌های هوشمند پردازش داده‌ها هستند که در ابتدا به وسیله پاپسکی و همکاران (۱۹۹۶)، شاپ و بوتن (۱۹۹۶) و تاماری و همکاران (۱۹۹۶) مورد استفاده قرار گرفته‌اند. فریت و همکاران (۲۰۰۳) در تخمین رطوبت خاک با استفاده از شبکه‌های عصبی مصنوعی گزارشی کردند که شبکه‌های عصبی مصنوعی ابزار مناسبی برای تخمین رطوبت خاک می‌باشند. پاپسکی و همکاران (۱۹۹۶) با مقایسه مدل‌های رگرسیونی و شبکه‌های عصبی مصنوعی در ۲۳۰ نمونه خاک گزارشی کردند که شبکه عصبی با دقت بیشتری ظرفیت نگهداری آب در خاک را برآورد می‌کند. شاپ و لچی (۱۹۹۸) نشان دادند که کیفیت پیش‌بینی توابع انتقالی توسط شبکه‌های عصبی وابسته به مجموعه داده‌هایی است که برای آموزش استفاده می‌شود و دقت برآوردها به‌طور مستقیم متأثر از کیفیت و نوع داده‌ها می‌باشد. تاواراکاوی و همکاران (۲۰۰۹) در اشتقاق توابع هیدرولیکی خاک با استفاده از روش‌های غیرپارامتری، بهبود چشم‌گیری را در تخمین‌های صورت گرفته مشاهده نمودند. یکی دیگر از روش‌های غیرپارامتریک که اخیراً در بسیاری از علوم از جمله کشاورزی (بنایان و هوگنبوم، ۲۰۰۹)، جنگل (جرتسن، ۲۰۰۷؛ لوپز و همکاران، ۲۰۰۱) و هیدرولوژی (کلارک و همکاران، ۲۰۰۴؛ یتس و همکاران، ۲۰۰۳) کاربرد چشم‌گیری یافته است مبتنی بر الگوشناسی و استفاده از اصل تشابه و نزدیکی داده‌ها است. چنین رویکردی در تخمین ویژگی‌های موردنظر، هنگامی که نوع رابطه بین داده‌های ورودی و خروجی روشن نباشد از سودمندی بیشتری برخوردار است (یاکویتز، ۱۹۹۳؛ لال و شارما، ۱۹۹۶). روش K نزدیک‌ترین همسایگی (K-NN)^۲ یکی از تکنیک‌های مبتنی بر چنین رویکردی است که در تخمین

1- Artificial Neural Network

2- K-Nearest Neighbor

برخی ویژگی‌های فیزیکی و شیمیایی خاک مورد استفاده قرار گرفته است. جلالی و همایی (۲۰۱۱) این روش را برای تخمین هدایت هیدرولیکی خاک با استفاده از برخی خصوصیات خاک مانند توزیع اندازه ذرات، هدایت الکتریکی اشباع، رطوبت اشباع، کربن آلی، مقدار مواد خشتی‌شونده، جرم ویژه حقیقی و ظاهری به‌کار بردند. آن‌ها گزارش کردند که این تکنیک در بیش‌تر موارد به شکل قابل‌قبولی توانایی تخمین کمیت موردنظر را دارد و می‌تواند جایگزین مناسبی برای اشتقاق توابع انتقالی خاک محسوب گردد. نمس و همکاران (۲۰۰۹) نیز با مقایسه کارایی توابع انتقالی پارامتریک و تکنیک K-NN در خاک‌های ایالت متحده، توانایی روش K-NN در تخمین توابع هیدرولیکی خاک را در مقیاس کل خاک‌های ایالت متحده بهتر از توابع پارامتریک ارزیابی نمودند. نمس و همکاران (۲۰۰۶) در برآورد نقاط پتانسیلی خاک (۳۳- و ۱۵۰۰- کیلوپاسکال) از روی توزیع اندازه ذرات، جرم ویژه ظاهری و درصد مواد آلی ۲۱۲۵ خاک از پایگاه داده امریکا (U.S. NRCS-SCS) توانایی دو روش ANN و K-NN را مورد بررسی قرار دادند و گزارش کردند که هر دو روش یاد شده از دقت یکسانی در تخمین و اشتقاق توابع توابع هیدرولیکی برخوردارند. آن‌ها همچنین کاربرد روش K-NN را در مواردی که توسعه دوباره توابع انتقالی نیاز نباشد به‌عنوان روش جایگزین مناسبی در اشتقاق این توابع پیشنهاد نمودند. در پژوهشی مشابه، حق‌وردی و همکاران (۲۰۱۰) نیز برای تخمین نقاط پتانسیلی آب در خاک (نقطه پژمردگی دائم و ظرفیت مزرعه) دو روش غیرپارامتریک ANN و K-NN را مورد مقایسه قرار دادند و گزارش کردند که روش K-NN از دقت بالاتری نسبت به ANN در خاک‌های شمال و شمال شرق ایران برخوردار است. از آن‌جا که تاکنون در ایران هیچ پژوهشی با هدف مقایسه عملکرد دو روش غیرپارامتریک K نزدیک‌ترین همسایگی (K-NN) و شبکه عصبی مصنوعی (ANN) بر روی CEC خاک‌های کشور گزارش نشده است، بنابراین در این پژوهش ضمن بررسی سودمندی کاربرد تکنیک K-NN در برآورد CEC برخی از خاک‌های ایران، دقت تخمین این کمیت با دو روش K-NN و ANN نیز مورد ارزیابی قرار خواهد گرفت.

مواد و روش‌ها

تعداد ۶۸۳ نمونه خاک از برخی استان‌های ایران شامل مناطقی از استان سمنان ۲۷۸ خاک، دشت قزوین ۱۶۷ خاک، ساوجبلاغ از استان البرز (۱۲۶ خاک) و ورامین (۱۱۲ خاک) از عمق ۲۵ سانتی‌متری سطحی، برداشت شد. از این تعداد ۱۲۰ عدد خاک به‌صورت تصادفی انتخاب و به‌عنوان خاک‌های

مورد آزمون (هدف) برای تخمین CEC در نظر گرفته شدند. همه نمونه‌ها پس از هوا خشک شدن و کوبیده شدن از الک ۲ میلی‌متری عبور داده شدند و برای انجام آزمون‌های فیزیکی و شیمیایی آماده شدند. توزیع اندازه ذرات به روش هیدرومتری (گی و بادر، ۱۹۸۶)، کربن آلی با روش اکسیداسیون تر والکلی - بلک (نلسون و سامرز، ۱۹۸۲) و اندازه‌گیری CEC نیز با روش استاندارد (چاپمن، ۱۹۶۵) انجام شد.

مدل شبکه عصبی مصنوعی (ANN): شبکه‌ای از عناصر پردازش ساده (نورون‌ها)، که می‌توانند رفتار پیچیده معینی را از ارتباط بین عناصر پردازش و پارامترهای عنصر نمایش دهند، اساس پردازش در این شبکه محسوب می‌شود. واحدهای محاسباتی شبکه (پرسپترون‌ها) می‌توانند به دو شکل تک‌لایه (SLP)^۱ و یا چندلایه (MLP)^۲ عمل نمایند. در این پژوهش از الگوریتم پرسپترون MLP با سه لایه برای برآورد داده‌های هدف استفاده شد. این نوع الگوریتم از یک لایه ورودی به منظور اعمال ورودی‌های مسئله برای آموزش شبکه با داده‌های توزیع اندازه ذرات، درصد کربن آلی و مقادیر ظرفیت تبادل کاتیونی این خاک‌ها، یک لایه پنهان با ۷ نورون و تابع عمل‌گر سیگموئیدی و یک لایه خروجی با تابع عمل‌گر خطی تشکیل شده است و قادر خواهد بود که برآوردی از ظرفیت تبادل کاتیونی خاک‌های مورد آزمون را با توجه به بانک داده‌های ورودی (مرجع) ارائه نمایند. برای سنجش عملکرد ANN در اشتقاق توابع و پیش‌بینی ظرفیت تبادل کاتیونی خاک‌ها، مقادیر آماره‌های ضریب همبستگی (r)، ریشه میانگین مربعات خطا (RMSE)، میانگین قدرمطلق خطا (MAE)، خطای ماکزیمم (ME)، میانگین باقی‌مانده‌ها (MR) و در نهایت کارایی مدل (EF) محاسبه شدند.

روش K نزدیک‌ترین همسایگی (K-NN): به‌طور کلی K-NN روشی برای طبقه‌بندی یک عنصر در یک مجموعه است که این کار را براساس نزدیک‌ترین خصوصیات سایر اعضای موجود در مجموعه (نمونه‌های آموزش‌دهنده) انجام می‌دهد. K-NN دارای یکی از ساده‌ترین الگوریتم‌های آموزش برای پیش‌بینی داده‌ها است و بر خلاف PTF های دیگر از هیچ تابع ریاضیاتی از پیش مشخص‌شده‌ای برای تخمین عناصر هدف استفاده نمی‌کند. الگوریتم K-NN به گروه موسوم به الگوریتم‌های یادگیرنده تنبل (lazy learning algorithms) تعلق دارد و داده‌ها را به‌گونه‌ای ذخیره می‌نماید که تا زمانی که دسترسی برای برآورد جدید نباشد، فرایند آموزش روی این داده‌ها انجام نخواهد شد (نمس و همکاران،

1- Single-Layer Perceptron

2- Multi-Layer Perceptron

۲۰۰۶؛ نمس و همکاران، ۲۰۰۸). در این تکنیک، تعداد بهینه‌ای از نمونه‌های موجود در یک مجموعه که دارای شبیه‌ترین ویژگی‌ها به نمونه هدف باشند (K) در نظر گرفته می‌شود و سپس جایابی و طبقه‌بندی عنصر هدف در این مجموعه از داده‌ها، با تعیین فاصله‌ها و سپس وزن‌دهی عناصر آموزش‌دهنده صورت می‌پذیرد. در یک بانک داده مرجع خاک، نزدیک‌ترین (مشابه‌ترین) خاک‌ها به خاک مورد آزمون انتخاب می‌شوند (K) و این خاک‌ها با توجه به میزان مشابهت به خاک مورد آزمون وزن‌دهی می‌شوند. به این ترتیب که در ابتدا فواصل اقلیدسی بین خاک مورد آزمون و خاک‌های مرجع محاسبه شده، سپس تعداد K عدد از نزدیک‌ترین خاک‌ها براساس فاصله به‌دست آمده، وزن‌دهی می‌شوند. در انتها با توجه وزن هر یک از خاک‌ها در مجموعه‌ای با K عدد خاک (نزدیک‌ترین همسایه)، برآوردی از داده‌های هدف که مورد آزمون واقع شده‌اند، صورت می‌پذیرد (نمس و همکاران، ۲۰۰۶).
با توجه به موارد بالا در این پژوهش به ترتیب زیر عمل گردید.
الف) محاسبه فاصله اقلیدسی بین نمونه هدف با هر یک از نمونه‌های خاک در بانک داده مرجع از رابطه زیر انجام شد:

$$d_i = \sqrt{\sum_{j=1}^x \Delta a_{ij}^2} \quad (1)$$

که در آن، d_i : نشان‌دهنده فاصله i امین خاک از بانک داده مرجع تا خاک هدف می‌باشد، Δa_{ij} : اختلاف فاصله

$$a_{ij(temp)} = \frac{[(a_{ij}) - \bar{a}_i]}{\sigma(a_j)} \quad (2)$$

که در آن، a_{ij} : نماینده j امین متغیر از i امین خاک است، \bar{a}_i و $\sigma(a_j)$ به ترتیب میانگین و انحراف معیار مقادیر مشاهده‌ای از j امین متغیر در بانک داده‌های مرجع است.

سپس مقادیر نهایی تبدیل شده j امین متغیر از i امین خاک ($a_{ij(trans)}$) با توجه به دامنه کم‌ترین تا بیش‌تر مقدار $a_{ij(temp)}$ از رابطه زیر به دست آمد و به جای مقادیر عددی واقعی به عنوان داده ورودی مورد استفاده قرار گرفت (نمس و همکاران، ۲۰۰۶).

$$a_{ij} = \frac{\{\max[range(a_{j=y(temp)}), \dots, range(a_{j=x(temp)})]\}}{range(a_{j(temp)})} \quad (3)$$

پس از نرمال کردن متغیرهای مستقل ورودی و محاسبه فاصله‌ها بین نمونه هدف و سایر نمونه‌های موجود در بانک داده مرجع، نمونه خاک‌های مورد مطالعه بر حسب مقدار محاسبه شده برای فاصله، به ترتیب صعودی مرتب شدند.

ب) یافتن تعداد بهینه خاک‌هایی که نزدیک‌ترین فاصله را با نمونه هدف دارند (مقدار K). به این منظور از روش ارزیابی تقاطعی (Cross Validation) استفاده شد. لال و شارما (۱۹۹۶) پیشنهاد کردند که اگر تعداد کل نمونه‌ها در بانک داده مرجع بیش از ۱۰۰ عدد باشد مقدار بهینه K برای رسیدن به کم‌ترین خطای برآوردها، تقریباً برابر ریشه دوم تعداد نمونه‌ها خواهد بود. نمس و همکاران (۲۰۰۶) در پژوهشی با اعمال اندازه‌های مختلف بانک داده مرجع از خصوصیات هیدرولیکی خاک نتیجه گرفتند که مقدار بهینه K به صورت تابعی از تعداد کل خاک‌های بانک مرجع (N) است. آن‌ها رابطه زیر را برای تعیین بهترین مقدار برای خاک پیشنهاد نمودند.

$$K = 0.765N^{0.493}, R^2 = 0.989 \quad (4)$$

در این پژوهش با توجه به تعداد نمونه‌های بانک مرجع مقدار خطای برآورد در دامنه‌ای از مقادیر مختلف K محاسبه شد و براساس روش ارزیابی تقاطعی، بهترین تعداد برای K ، زمانی که خطای تخمین‌ها (RMSE)^۱ به کم‌ترین مقدار خود رسید تعیین گردید.

1- Root Mean Square Error

ج) وزن‌دهی نمونه خاک‌های آموزش‌دهنده (K) براساس فاصله‌های به‌دست آمده بین خاک هدف و هر یک از خاک‌های تعیین شده در K طبق رابطه‌های زیر صورت گرفت.

$$w_i = \frac{d_{i(rel)}}{\sum_{i=1}^K d_{i(rel)}} \quad (5)$$

که در آن، w_i : وزن i امین نمونه از خاک‌های آموزش‌دهنده (K) است و $d_{i(rel)}$ فاصله نسبی هدف و نمونه i می‌باشد که از رابطه زیر به‌دست آمده است.

$$d_{i(rel)} = \left(\frac{\sum_{i=1}^K d_i}{d_i} \right)^p \quad (6)$$

نمس و همکاران (۲۰۰۶) کم‌ترین مقدار خطای تخمین در توابع هیدرولیکی خاک را در p های نزدیک به ۱ به‌دست آوردند و گزارش نمودند که حساسیت روش K-NN به مقادیر مختلف p بسیار اندک می‌باشد. در این پژوهش علاوه بر بهینه‌سازی تعداد نمونه خاک‌های نزدیک‌ترین همسایه (K)، بهترین مقدار p نیز از تکنیک ارزیابی تقاطعی به‌دست آمد.

د) برآورد مقادیر CEC در نمونه‌های هدف از طریق محاسبه میانگین وزنی CEC در نمونه خاک‌های K با استفاده از رابطه زیر صورت گرفت.

$$P_i = \sum_{i=1}^K w_i O_i \quad (7)$$

که در آن، P_i : مقدار CEC پیش‌بینی شده خاک هدف، w_i : وزن مربوط به i امین خاک در K و O_i مقدار CEC اندازه‌گیری شده i امین خاک در K است.

ه) در آخرین مرحله مقادیر آماره‌های لازم برای مقایسه و ارزیابی قابلیت دو روش ANN و K-NN در پیش‌بینی CEC خاک‌های مختلف از رابطه‌های زیر به‌دست آمد.

$$r = \frac{n \left(\sum_{i=1}^n (P_i O_i) \right) - \left(\sum_{i=1}^n P_i \right) \left(\sum_{i=1}^n O_i \right)}{\sqrt{\left[n \sum_{i=1}^n (P_i)^2 - \left(\sum_{i=1}^n P_i \right)^2 \right] \left[n \sum_{i=1}^n (O_i)^2 - \left(\sum_{i=1}^n O_i \right)^2 \right]}} \quad (8)$$

که در آن، r : ضریب همبستگی پیرسون، P_i و O_i : به ترتیب مقادیر CEC پیش‌بینی شده و مشاهده‌ای و n : تعداد نمونه‌های مورد آزمون است. این ضریب که بین -1 تا 1 تغییر می‌کند بیان‌کننده میزان هم‌روندی داده‌هاست. هرچه این مقدار به 1 و یا -1 نزدیک‌تر باشد، نشانه آن است که روند تغییرات مقادیر پیش‌بینی شده و مشاهده‌ای یکنواخت‌تر می‌باشد. به منظور مقایسه عملکرد دو روش ANN و K-NN آماره‌های ریشه میانگین مربعات خطا (RMSE)، میانگین قدرمطلق خطا (MAE)^۱، بیش‌ترین خطا (MaxE)^۲، میانگین باقی‌مانده (ME)^۳ و کارایی مدل (EF)^۴ با استفاده از رابطه‌های زیر به دست آمدند.

$$RMSE = \sqrt{\frac{\sum_{i=1}^n (P_i - O_i)^2}{n}} \quad (9)$$

$$MAE = \sqrt[n]{\sum_{i=1}^n |P_i - O_i|} \quad (10)$$

$$ME = \max |P_i - O_i| \quad (11)$$

$$MR = \sqrt[n]{\sum_{i=1}^n (P_i - O_i)} \quad (12)$$

$$EF = \frac{\sum_{i=1}^n (o_i - \bar{o})^2 - \sum_{i=1}^n (p_i - o_i)^2}{\sum_{i=1}^n (o_i - \bar{o})^2} \times 100 \quad (13)$$

که در آن‌ها، \bar{o} : میانگین مقادیر مشاهده‌ای و CEC می‌باشد.

در این پژوهش برنامه کامپیوتری در محیط برنامه‌نویسی Matlab (7.10) تهیه شد (Matlab 7.1, The Math works Inc., Natick, MA) و تمامی مراحل (الف) تا (ه) با استفاده از این برنامه به انجام رسید.

-
- 1- Mean Absolute Error
 - 2- Maximum Error
 - 3- Mean Error
 - 4- Efficiency of Model

نتایج و بحث

واکنش خاک (pH) در نمونه‌های مورد مطالعه بین ۷/۶-۸/۳ قرار داشت. محدوده هدایت الکتریکی (EC) نمونه‌ها ۴-۱۱ دسی‌زیمنس بر متر و نسبت جذب سدیم (SAR) آن‌ها نیز کم‌تر از ۱۳ بود که نشان‌دهنده شوری متوسط تا زیاد خاک‌های مورد مطالعه بود. جدول ۱ خلاصه‌ای از ویژگی‌های اندازه‌گیری شده خاک‌های هدف و بانک مرجع را نشان می‌دهد. همان‌گونه که مشاهده می‌شود، خاک‌های مورد مطالعه در خصوصیات مانند توزیع اندازه ذرات و کربن آلی که از مهم‌ترین عوامل موثر بر CEC خاک‌ها هستند، دامنه به‌نسبت وسیعی را نشان می‌دهند. این امر را می‌توان از دلایل اصلی تنوع مقادیر CEC اندازه‌گیری شده در این پژوهش نیز به حساب آورد. جدول ۱، میانگین و انحراف معیار تقریباً یکسانی را بین خصوصیات خاک‌های مورد آزمون و بانک داده مرجع نشان می‌دهد. این امر می‌تواند به بهبود تخمین‌های CEC با هر دو روش غیرپارامتریک مورد مقایسه در این پژوهش کمک نماید. از آن‌جا که الگوریتم روش ANN از تمام داده‌های بانک مرجع و ارتباط بین آن‌ها برای آموزش استفاده می‌نماید، همگنی داده‌ها می‌تواند کمک قابل توجهی به بهبود نتایج برآوردها با این روش نماید این شرایط برای روش K-NN که از تعداد محدودتری از داده‌ها که کم‌ترین فاصله را با داده‌های هدف دارند (K)، به مراتب منافع کمتری را به‌دنبال خواهد داشت. حق‌وردی و همکاران (۲۰۱۰) در تخمین توابع هیدرولیکی خاک با مقایسه اثر کیفیت و نوع داده‌ها بر عملکرد دو روش ANN و K-NN گزارش کردند که حساسیت روش ANN به ناهمگنی داده‌ها بیش‌تر از روش K-NN بوده به‌طوری‌که همگن نبودن داده‌ها در این پژوهش باعث افزایش ۱۰۰ درصدی خطای برآوردها در این روش شد.

بهینه‌سازی هم‌زمان K و P با استفاده از روش ارزیابی تقاطعی به گونه‌ای انجام شد که به‌ازای هر مقدار P از ۲-۰/۲ (با فاصله‌های ۰/۲) K های مختلف از ۱ تا ۴۰ در نظر گرفته شد و برای هر ترکیب P و K برای هر یک از ۱۲۰ عدد خاک‌های مورد آزمون، برنامه کامپیوتری تهیه شده، یک بار اجرا گردید و مقادیر RMSE به‌دست آمد. شکل ۱ نمودار سه‌بعدی تغییرات RMSE را با تغییرات هم‌زمان K و P نشان می‌دهد. نتایج این بررسی نشان داد که کم‌ترین خطای تخمین CEC در P=۱ و K=۱۴ وجود دارد. نمس و همکاران (۲۰۰۶) نتیجه گرفتند که P می‌تواند اثر منفی کوچک بودن بانک داده مرجع را در عملکرد الگوریتم بهبود بخشد. در حقیقت مقادیر غیر از ۱ برای P، با تأثیر بر وزن‌دهی نقاط سبب تعدیل اثرات منفی کوچک بودن بانک داده مرجع می‌شود. به‌طوری‌که در مقادیر کوچک P ($P < 1$)، اثر وزن نقاط نزدیک به داده هدف کاهش و با افزایش P به بزرگ‌تر از ۱، این اثر افزایش

می‌یابد. در این پژوهش با توجه به نوع و تعداد خاک‌ها در بانک داده مرجع، کم‌ترین RMSE در $P=1$ به‌دست آمد (شکل ۳). بر این اساس رابطه‌های ۵ و ۶ در $P=1$ به‌صورت ساده زیر بازنویسی شدند و بهترین وزن‌ها برای خاک‌های K از رابطه زیر به‌دست آمد.

$$w_i = \frac{1/d_i}{\sum_{i=1}^K 1/d_i} \quad (14)$$

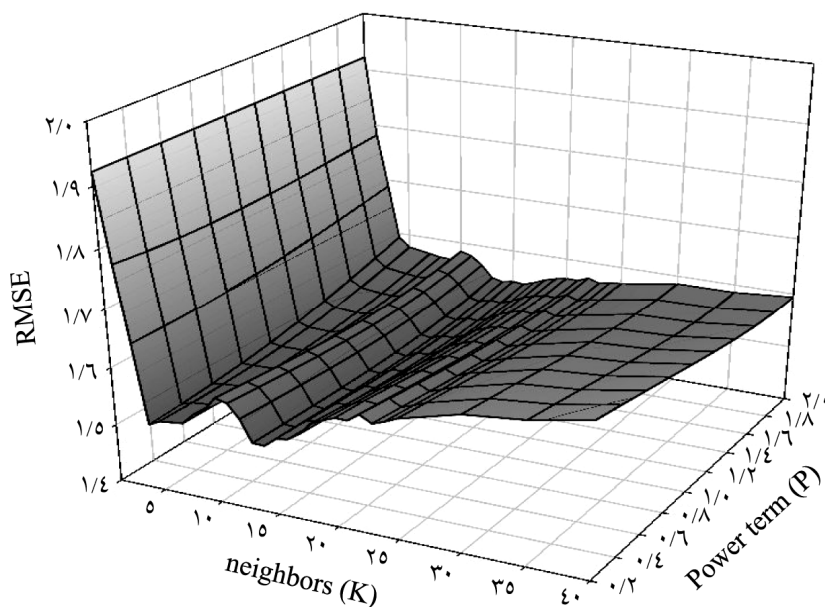
جدول ۱- خلاصه آماری برخی از ویژگی‌های خاک‌های بانک مرجع و هدف.

ضریب تغییرات (درصد)	میان	انحراف معیار	میانگین	مینیمم	ماکزیمم	واحد	خصوصیت	
۵۲/۸	۱۹/۰۰	۱۲/۶۴	۲۳/۹۲	۴/۷۲	۶۸/۰۰		رس	بانک داده مرجع (N=۵۶۳)
۲۷/۶	۴۲/۰۰	۱۱/۱۶	۴۰/۴۵	۴/۰۰	۷۲/۰۰	(درصد)	سیلت	
۳۷/۶	۳۶/۰۰	۱۳/۴۸	۳۵/۸۶	۵/۳۸	۹۱/۲۸		شن	
۷۰/۷	۰/۵۹	۰/۵۸	۰/۸۲	۰/۰۱	۳/۶۹		کربن آلی	
۳۵/۹	۹/۷۹	۴/۴۲	۱۲/۳۰	۶/۵۰	۲۹/۴۰	$\text{Cmol}^+ \text{kg}^{-1}$	ظرفیت تبادل کاتیونی	
۵۹/۸	۱۸/۰۰	۱۳/۸۹	۲۳/۲۱	۵/۹۲	۶۴/۲۵		رس	داده‌های هدف (N=۱۲۰)
۲۶/۵	۴۳/۵۱	۱۱/۰۲	۴۱/۶۳	۶/۰۰	۷۳/۰۰	(درصد)	سیلت	
۳۷/۶	۳۶/۰۰	۱۳/۵۶	۳۶/۰۴	۴/۷۲	۷۶/۵۶		شن	
۶۹/۸	۰/۵۹	۰/۶۰	۰/۸۶	۰/۰۲	۲/۶۱		کربن آلی	
۳۶/۹	۹/۵۲	۴/۴۵	۱۲/۰۶	۶/۹۰	۲۳/۲۶	$\text{Cmol}^+ \text{kg}^{-1}$	ظرفیت تبادل کاتیونی	

N: تعداد داده‌ها است.

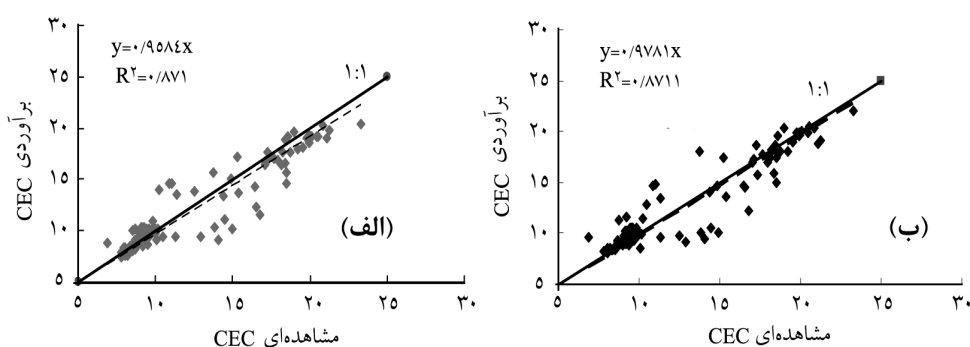
شکل ۲ مقادیر به‌دست آمده RMSE را در $P=1$ در مقابل K های مختلف نشان می‌دهد، که بیش‌ترین مقدار خطا مطابق انتظار در حالتی که فقط یک خاک به‌عنوان نزدیک‌ترین همسایه برای تخمین‌ها استفاده شد ($K=1$) به‌دست آمد و کم‌ترین آن در حالت بهینه در $K=14$ وجود داشت و با

افزایش مقدار K به بیش از ۱۴ خطای برآوردها دوباره افزایش یافت. در این پژوهش مقدار K به دست آمده از طریق ارزیابی تقاطعی هم‌خوانی بسیار خوبی با نتایج نمس و همکاران (۲۰۰۶) نشان می‌دهد (رابطه ۴) و با آنچه لال و شارما (۱۹۹۶) توصیه نمودند، مغایرت داشت ($K \approx 24 = 563^{1/5}$). شکل ۳ نیز مقادیر خطای برآورد را در مقابل P های مختلف ارایه می‌کند. با توجه به تغییرات بسیار اندک RMSE (در حدود ۰/۰۳)، این شکل بیان می‌کند که روش K -NN نسبت به تغییرات P حساسیت پایینی را در محدوده خاک‌های مورد مطالعه نشان می‌دهد و بنابراین انتخاب هر یک از مقادیر P در محدوده اعداد نزدیک به ۱ تفاوت قابل‌توجهی را در نتایج تخمین‌ها ایجاد نمی‌کند. البته تعداد خاک‌های موجود در بانک داده مرجع عامل مهمی در تعیین P بهینه می‌باشد. به‌عبارت دیگر استفاده از رابطه ساده (۱۴) به‌منظور وزن‌دهی خاک‌های K در بانک داده‌های مشابه می‌تواند نتایج قابل‌قبولی را ارایه نماید. نمس و همکاران (۲۰۰۶) در برآورد توابع هیدرولیکی با ۵ اندازه مختلف بانک داده مرجع مقادیر بین ۰/۹۵-۱/۱ را برای P بهینه به‌دست آوردند.



شکل ۱- رابطه مجذور میانگین مربعات خطا (RMSE)، با تغییرات هم‌زمان K و P در یک نمودار سه‌بعدی.

CEC پیش‌بینی شده در مقابل مقادیر مشاهده‌ای این کمیت در روش شبکه عصبی مصنوعی (الف) و K نزدیک‌ترین همسایگی (ب) ترسیم شده است. همان‌طور که مشاهده می‌شود شیب معادله‌های خطوط برازش داده شده بر این نقاط تفاوت بسیار اندکی را با شیب خط ۱:۱ ترسیمی نشان می‌دهد. در این جا ضریب تبیین بالا (۰/۸۷) با شیب نزدیک به ۱، بی‌تردید بیانگر دقت بالا و یکسان هر دو روش در تخمین CEC خاک‌های مورد مطالعه می‌باشد. ضریب همبستگی بالا و بسیار معنی‌دار ۰/۹۴ (در سطح احتمال ۹۹ درصد) نشان‌دهنده روند یکسان تغییرات مقادیر برآوردی و مشاهده‌ای در روش K-NN بود که از این نظر هم تفاوتی با روش ANN نشان نداد. به‌منظور مقایسه عملکرد دو روش ANN و K-NN آماره‌های RMSE، MAE، ME، MR و EF را نیز مورد بررسی قرار گرفت. جدول ۲ این مقادیر را برای دو روش مورد آزمون نشان می‌دهد. نتایج نشان داد که در روش K-NN با $K=14$ و $P=1$ در بدترین حالت، CEC را با $cmol^+/kg$ ۸۱/۴ خطا (ME)، پیش‌بینی نمود. این شرایط در روش ANN با تخمین ضعیف‌تر و میزان خطای $cmol^+/kg$ ۲۶/۵ همراه بود. مقادیر کوچک آماره‌های RMSE و MAE در هر دو روش K-NN و ANN نشان داد که این دو روش در برآورد CEC خاک‌های مورد آزمون از دقت بالا و یکسانی برخوردارند. اگرچه این مقدار در روش K-NN از نظر عددی اندکی بهتر از روش ANN بود. مقادیر مثبت آماره MR برای این دو روش نیز نشان داد که هر دوی آن‌ها متمایل به برآورد کم‌تر مقدار پیش‌بینی شده نسبت به داده‌های مشاهده‌ای (Under estimation) می‌باشند. همچنین مقایسه مقادیر EF نشان داد که هر دو روش دارای کارایی بالایی ($EF=88$ درصد) در برآورد ظرفیت تبادل کاتیونی از روی توزیع اندازه، درصد کربن آلی خاک می‌باشند و هیچ اختلافی از این نظر بین دو روش مشاهده نمی‌شود.



شکل ۴- رابطه بین مقادیر پیش‌بینی و مشاهده‌ای CEC با استفاده از دو روش ANN (الف) و K-NN (ب).

جدول ۲- آماره‌های مقایسه عملکرد دو روش K نزدیک‌ترین همسایگی (K-NN) و شبکه عصبی مصنوعی (ANN).

روش	r	RMSE	MAE	MaxE	ME	EF
K-NN	۰/۹۴	۱/۵۱	۰/۹۲	۴/۸۱	۰/۱	۸۸
ANN	۰/۹۴	۱/۵۳	۱/۰۰	۵/۲۶	۰/۱۸	۸۸

ریشه میانگین مربعات خطا (RMSE)، میانگین قدرمطلق خطا (MAE)، بیش‌ترین خطا (MaxE)، میانگین خطا (ME) بر حسب cmol^+/kg و کارایی مدل (EF) بر حسب درصد می‌باشد.

نتیجه‌گیری

از آن‌جا که تاکنون در ایران هیچ پژوهشی با هدف مقایسه عملکرد دو روش غیرپارامتریک K نزدیک‌ترین همسایگی (K-NN) و شبکه عصبی مصنوعی (ANN) بر روی CEC خاک‌های کشور گزارش نشده است، ضرورت چنین بررسی مورد توجه نویسندگان قرار گرفت. بر این اساس در این پژوهش قابلیت دو روش نام برده در برآورد CEC خاک از روی ویژگی‌های زودیافتی مانند بافت و درصد کربن آلی خاک مورد ارزیابی و مقایسه قرار گرفت. بررسی ضرایب همبستگی در هر دو روش مورد مقایسه، هم‌روندی بالا و بسیار معنی‌داری ($P < 0/01$) را بین مقادیر مشاهده‌ای و برآوردی CEC خاک‌ها نشان داد. در مدل K-NN که با استفاده از تشابه ۱۴ عدد از نزدیک‌ترین خاک‌های بانک مرجع به هدف، مقدار CEC را پیش‌بینی می‌کرد، بین مقادیر CEC برآوردی و اندازه‌گیری شده (نسبت به خط ۱:۱) ضریب تبیین (R^2) ۰/۸۷ مشاهده شد که به روشنی توانایی و دقت بسیار بالای این روش را در تخمین CEC خاک از روی ویژگی‌های یاد شده نشان می‌دهد. این دقت در تخمین CEC خاک‌ها در روش نام برده، با مدل ANN که از ۵۶۳ داده بانک مرجع برای آموزش استفاده می‌نمود، یکسان بود، به طوری که مقادیر RMSE محاسبه شده برای روش‌های K-NN و ANN به ترتیب ۱/۵۱ و ۱/۵۳ (cmol^+/kg) به دست آمد. مقدار مثبت میانگین خطاهای (ME) محاسبه شده نیز نشان داد که به طور کلی هر دو روش، مقدار CEC را کم‌تر از مقدار مشاهده‌ای این کمیت تخمین زده‌اند. کارایی (EF) ۸۸ درصدی برای هر دو روش، قابلیت بالا و یکسان آن‌ها را در تخمین CEC خاک‌ها نشان می‌دهد. بر این اساس و با توجه به کاربرد آسان روش K-NN در پیش‌بینی CEC خاک‌ها به خصوص زمانی که داده‌های موجود همگن نباشند و یا تعداد این داده‌ها برای آموزش نوروهای شبکه عصبی کافی نباشد، استفاده از روش K-NN توصیه می‌شود.

منابع

1. Amini, M., Abbaspour, K.C., Khademi, H., Fathianpour, N., Afyuni, M., and Schulin, R. 2005. Neural Network Models to Predict cation exchange capacity in arid regions of Iran. *Eur. J. Soil Sci.* 56: 551-559.
2. Bannayan, M., and Hoogenboom, G. 2009. Using pattern recognition for estimating cultivar coefficients of a crop simulation model. *Field Crop. Re.* 111: 290-302.
3. Bell, M.A., and Van Kullen, H. 1995. Soil pedotransfer function for four Mexican soils. *Soil. Sci. Soc. Am. J.* 59: 865-871.
4. Carrow, R.N., Waddington, D.V., and Rieke, P.E. 2001. Turfgrass soil fertility and chemical problems. John Wiley and Sons, New York.
5. Chapman, H.D. 1965. Cation exchange capacity. *Methods of Soil Analysis. Part 2. Chemical and microbiological properties.* ASA-CSSA-SSSA Publisher, Madison, Wisconsin, USA.
6. Clark, M.P., Gangopadhyay, S., Brandon, D., Werner, K., Hay, L., Rajagopalan, B., and Yates, D. 2004. A resampling procedure for generating conditioned daily weather sequences. *Water Resour. Res.* 40. W04304 doi. 10.1029/2003WR002747.
7. Frate, F.D., Ferrazoli, P., and Schiavon, G. 2003. Retrieving soil moisture and agricultural variables by microwave radiometry using neural network. *Remote Sens. Environ.* 84: 174-183.
8. Gee, G.W., and Bauder, J.W. 1986. Particle-size analysis. *Methods of Soil Analysis Part 1. Physical and mineralogical methods,* ASA-CSSA-SSSA Publisher, Madison, Wisconsin, USA.
9. Gjertsen, A.K. 2007. Accuracy of forest mapping based on Landsat TM data and a kNN-based method. *Remote Sens. Environ.* 110: 420-430.
10. Haghverdi, A., Ghahraman, B., Khoshnood Yazdi, A.A., and Arabi, Z. 2010. Estimating of water content in FC and PWP in north and north east of Iran's soil samples using k-nearest neighbor and artificial neural networks. *J. Water Soil.* 24: 4. 804-814. (In Persian)
11. Jalali, V.R., and Homae, M. 2011. A nonparametric model by using k-nearest neighbor technique for predicting soil saturated hydraulic conductivity. *J. Water Soil.* 25: 2. 347-355. (In Persian)
12. Karimian, N.A. 1998. Clay and organic matter contribution in cation exchange capacity of calcareous soils from Fars province. *Proceedings of 5th Iranian soil congress, Karaj, Iran.* (In Persian)
13. Krogh, L., Madsen, H.B., and Greve, M.H. 2000. Cation exchange capacity pedotransfer functions for Danish soils. *Acta Agric. Scand., Sect. B, Soil and Plant Sci.* 50: 1-12.
14. Lall, U., and Sharma, A. 1996. A nearest-neighbor bootstrap for resampling hydrologic time series. *Water Resour. Res.* 32: 679-693.

15. Lopez, H.F., Ek, A.R., and Bauer, M.E. 2001. Estimation and mapping of forest stand density, volume, and cover type using the k-nearest neighbors method. *Remote Sens. Environ.* 77: 251-274.
16. Mc Bratney, A.B., Minasny, B., Cattle, S.R., and Vervoort, R.W. 2002. From pedotransfer functions to soil inference systems. *Geoderma*. 109: 41-73.
17. Memarian Fard, M., and Beigi Harchagani, H. 2009. Comparison of artificial neural network and regression pedotransfer functions models for prediction of network and regression pedotransfer functions models for prediction of J. *Water Soil*. 23: 4. 90-99. (In Persian)
18. Nelson, D.W., and Sommers, L.E. 1982. Total carbon, organic carbon, and organic matter. *Methods of Soil Analysis. Part 2. Chemical and microbiological properties*. ASA-CSSA-SSSA Publisher, Madison. Wisconsin, USA.
19. Nemes, A., Rawls, W.J., and Pachepsky, Ya.A. 2006. Use of the nonparametric nearest neighbor approach to estimate soil hydraulic properties. *Soil Sci. Soc. Am. J.* 70: 327-336.
20. Nemes, A., Roberts, R.T., Rawls, W.J., Pachepsky, Ya.A., and Van Genuchten, M.Th. 2008. Software to estimate -33 and -1500 kPa soil water retention using the non-parametric k-Nearest Neighbor technique. *Environm Modell. Softw.* 23: 254-255.
21. Nemes, A., Timlin, D.J., Pachepsky, Ya.A., and Rawls, W.J. 2009. Evaluation of the Rawls et al. (1982) pedotransfer functions for their applicability at the U.S. national scale. *Soil. Sci. Soc. Am. J.* 73: 1638-1645.
22. Pachepsky, Ya.A., Timlin, D., and Varallyay, G. 1996. Artificial neural networks to estimate soil water retention from easily measurable data. *Soil Sci. Soc. Am. J.* 6: 727-733.
23. Schaap, M.G., and Bouten, W. 1996. Modeling water retention curves of sandy soils using neural networks. *Water Resour. Res.* 32: 3033-3040.
24. Schaap, M.G., and Leij, F.L. 1998. Database-related accuracy and uncertainty of pedotransfer functions. *Soil Science*, 10: 765-779.
25. Seybold, C.A., Grossman, R.B., and Reinsch, T.G. 2005. Predicting cation exchange capacity for soil survey using linear models. *Soil Sci. Soc. Am. J.* 69: 856-863.
26. Tamari, S., Wosten, J.H.M., and Ruize, J.C. 1996. Testing an artificial neural networks for predicting soil hydraulic conductivity. *Soil. Sci. Soc. Am. J.* 60: 1732-1741.
27. Twarakavi, J., and Schaap, M.G. 2009. Development of pedotransfer functions for estimation of soil hydraulic parameters using support vector machines. *Soil. Sci. Soc. Am J.* 73: 1443-1452.
28. Yakowitz, S. 1993. Nearest-neighbor estimation for null-recurrent Markov time series. *Stoch. Proc. Appl.* 48: 311-318.
29. Yates, D., Gangopadhyay, S., Rajagopalan, B., and Strzepek, K. 2003. A technique for generating regional climate scenarios using a nearest-neighbor algorithm. *Water Resour. Res.* 39:1199 doi 10.1029/2002WR001769.



Comparison of K-nearest neighbor and artificial neural network methods for predicting cation exchange capacity of soil

***A.A. Zolfaghari¹, M.T. Tirgar Soltani², T. Afshari Badrloo³
and F. Sarmadian⁴**

¹Ph.D. Student, Dept. of Soil Science Engineering, University of Tehran,

²M.Sc., Dept. of Soil Science Engineering, University of Tehran,

³M.Sc. Student, Dept. of Soil Science Engineering, Shahrood University of Technology,

⁴Associate Prof., Dept. of Soil Science Engineering, University of Tehran

Received: 05/29/2012; Accepted: 12/15/2012

Abstract

Cation exchange capacity (CEC) measurement is a very expensive and time-consuming method in large scale assessments. It can be an appropriate approach to predict CEC from readily available properties via developing nonparametric models. In the present study, a nonparametric technique has been used for estimating CEC and compared with the most common nonparametric models which is based on artificial neural networks (ANN). 683 soils were selected from central Iran that 120 of them were used as target data and the others (563) were the reference data set. The parameters, clay, silt, sand and organic carbon content were the input independent variables (readily available properties) and the CEC was as an output dependent variable in this work. The results showed that the maximum error (MaxE) in K-NN and ANN techniques were 4.81 and 5.26 cmol⁺/kg, respectively. Root mean squared error (RMSE) for the K-NN and ANN were 1.51 and 1.53, respectively. This indicated that both methods are able to properly and equally predict CEC. The positive values of mean error (ME) showed that both models tended to underestimate CEC values in samples. The results analysis also showed that the efficiency of models (EF=0.88) were high by the estimation of CEC values in target soils.

Keywords: Nonparametric models, Cation exchange capacity, K-nearest neighbor, Artificial neural network model

* Corresponding Authors; Email: azolfaghari@ut.ac.ir