

## استفاده از مدل میانگین (مدل تجمیعی) در نقشه برداری رقومی خاک

شاهرخ فاتحی<sup>۱</sup>

استادیار بخش تحقیقات خاک و آب، مرکز تحقیقات و آموزش کشاورزی و منابع طبیعی استان کرمانشاه، سازمان تحقیقات و آموزش ترویج کشاورزی، کرمانشاه

### چکیده

یکی از راه‌های بهبود کارایی مدل‌سازی رقومی ویژگی‌های خاک که در نقشه‌برداری رقومی خاک مورد توجه قرار نگرفته‌است استفاده از مدل میانگین یا مدل تجمیعی می‌باشد. پژوهش حاضر، به منظور مقایسه مدل‌های پیش‌بینی منفرد مثل مدل رگرسیون خطی چندمتغیره، کوپست و جنگل تصادفی با روش مدل میانگین از نوع ترکیب خطی تعمیم‌یافته انجام شده‌است. ابتدا در زیرحوضه آبخیز مرک واقع در استان کرمانشاه به وسعت ۲۴۰۰۰ هکتار، ۳۲۰ نمونه خاک سطحی به روش تصادفی جمع‌آوری گردید و پس از تعیین کربن آلی خاک، داده‌ها به دو دسته آموزشی (۲۴۰ نمونه) و اعتبارسنجی (۸۰ نمونه) تفکیک شدند. متغیرهای کمکی، مشتقات مدل رقومی ارتفاع و شاخص پوشش گیاهی بودند. ریشه میانگین مربعات خطا در مدل میانگین (۰/۵۵۸) پایین‌تر از مدل رگرسیون خطی چندمتغیره (۰/۵۸۶)، کوپست (۰/۵۸۱) و جنگل تصادفی (۰/۵۷۵) به‌دست آمد. به‌طور کلی می‌توان بیان نمود استفاده از مدل میانگین در نقشه‌برداری رقومی خاک باعث افزایش صحت و دقت نقشه‌های پیش‌بینی می‌شود.

واژه‌های کلیدی: مدل‌سازی رقومی خاک، داده کاوی، مدل تجمیعی

### مقدمه

نقشه‌برداری رقومی خاک را مک برتنی و همکاران (۲۰۰۳) با استفاده از معادله ینی (۱۹۴۱) که عوامل خاکسازي را توصیف می‌کند، صورت بندی نمودند. که اصطلاحاً مدل اسکورپن نامیده می‌شود و به صورت معادله زیر بیان می‌گردد:

$$S_c[x.y.\sim t].S_p[x.y.\sim t] = f(s[x.y.\sim t].c[x.y.\sim t].o[x.y.\sim t].r[x.y.\sim t].p[x.y.\sim t].a[x.y.\sim t].n) \quad (1)$$

در این معادله،  $S_c$  = کلاس خاک،  $S_p$  = ویژگی‌های خاک،  $s$  = خاک‌ها (مشخصاتی از خاک در یک نقطه مشاهداتی)،  $c$  = اقلیم (خصوصیات اقلیمی محیط در یک نقطه مشاهداتی)،  $o$  = موجودات زنده (پوشش گیاهی، جانوران، یا فعالیت بشر)،  $r$  = توپوگرافی (مشخصات سیمای اراضی)،  $p$  = مواد مادری (سنگ شناسی)،  $a$  = سن (عامل زمان)،  $n$  = مکان (موقعیت مکانی)،  $t$  = زمان (اینجا منظور یک زمان تقریبی است)،  $x,y$  = مختصات مکانی و  $f$  = تابع یا تابع پیش‌بینی مکانی خاک می‌باشد. تابع پیش‌بینی مکانی در واقع روش‌های مختلف داده‌کاوی، یا زمین‌آماري و آماری هستند که برای استخراج روابط بین ویژگی‌های خاک و متغیرهای محیطی به کار می‌روند.

توابع پیش‌بینی مکانی یا همان روش‌های مدل‌سازی خاک بسته به ماهیت فرایند (خطی بودن و یا غیرخطی بودن روابط ویژگی‌های خاک و متغیرهای محیطی) و ویژگی داده‌ها (قابلیت دسترسی داده‌های کمی و کیفی) به کار می‌روند. رگرسیون خطی چندمتغیره، مدل‌های خطی تعمیم یافته (McCullagh and Nelder, 1989)، درختان طبقه بندی و رگرسیون (Breiman et al., 1984)، شبکه عصبی مصنوعی (Anderson, 1995) ماشین بردار پشتیبان (Cortes and Vapnik, 1995) و الگوریتم ژنتیک (Koza, 1992) مدل‌های مرسوم و رایجی هستند که در نقشه‌برداری رقومی خاک مورد استفاده قرار می‌گیرند. برخی از این مدل‌ها علاوه بر پیچیدگی، خروجی‌هایی ارائه می‌دهند که تفسیر آنها کم و بیش مشکل است (مانند شبکه عصبی). امامدل‌های

دیگر، دارای پیچیدگی متفاوتی هستند مانند مدل کوبیست و جنگل تصادفی از خانواده مدل درختان طبقه بندی و رگرسیون، که به طور خودکار بر اساس عملکرد زیرمدل‌ها، متغیرهای پیش‌بینی کننده در آنها انتخاب می‌شود. در پژوهش‌های نقشه برداری رقومی خاک معمولاً از چند مدل آماری، داده کاوی یا زمین آماری برای تهیه نقشه پیش‌بینی ویژگی‌های خاک استفاده می‌شود و در نهایت یک مدل بر اساس نتایج اعتبار سنجی به عنوان مدل برتر و دقیق تر انتخاب می‌گردد. حال این سؤال مطرح می‌شود آیا استفاده از ترکیبی از مدل‌های مختلف نسبت به مدل‌های منفرد دقیق تر نخواهد بود؟

چنین وضعیتی در رشته‌های پژوهشی هیدرولوژی و اتمسفر مشاهده می‌شود در این منظومه‌های علمی، پیش‌بینی‌های چندگانه از یک فرایند توسط تعدادی از مدل‌های پیش‌بینی کننده به دست می‌آید. هر مدل سهم در پیش‌بینی، نقاط قوت و ضعف خاص خود را دارد. ترکیب این مدل‌ها مسلماً بهتر از هر کدام از مدل پیش‌بینی به طور منفرد خواهد بود (Diks and Vrugt, 2010). در زمینه نقشه برداری رقومی خاک پادریان و همکاران (۲۰۱۴) و ملانو و همکاران (۲۰۱۴) از روش‌های مختلف مدل میانگین (مدل تجمیعی)<sup>۱</sup> مانند مدل میانگین با اوزان معادل<sup>۲</sup>، روش مدل میانگین واریانس وزن دار (روش باتز-گرانگر)<sup>۳</sup>، میانگین روش مدل میانگین گرانگر-رامانتان<sup>۴</sup> و روش مدل میانگین بیزی<sup>۵</sup> برای پیش‌بینی ویژگی‌های خاک استفاده کرده‌اند. هدف از پژوهش حاضر استفاده از مدل میانگین یا مدل تجمیعی برای پیش‌بینی کربن آلی و مقایسه آن با مدل‌های منفرد رگرسیون چند متغیره، کوبیست و جنگل تصادفی است.

## مواد و روش‌ها

اراضی مورد مطالعه با وسعتی حدود ۲۴۰۰۰ هکتار در بیست کیلومتری جنوب شرقی شهر کرمانشاه و در بین مختصات جغرافیائی ۴۷ درجه و ۴ دقیقه تا ۴۷ درجه و ۲۲ دقیقه طول شرقی و ۳۴ درجه تا ۳۴ درجه و ۹ دقیقه عرض شمالی قرار گرفته است. مطابق با آمار هواشناسی، متوسط درجه حرارت سالیانه هوا ۱۳/۲ درجه سانتیگراد و میانگین بارندگی سالیانه ۴۸۱/۴ میلیمتر است. بر اساس اقلیم نمای آمبروزه، اقلیم این ناحیه نیمه خشک و سرد می‌باشد. رژیم رطوبتی خاک زیریک (xeric) و رژیم حرارتی خاک ترمیک (thermic) است. سنگ آهک، دولومیت، مارن، سنگ رس و ماسه سنگ، سنگ‌های اصلی سازندهای زمین شناسی محدوده مطالعاتی هستند. خاک‌های اینسپتی سولز، انتی سولز و ورتی سولز مهم‌ترین رده‌های خاک شناسایی شده در این ناحیه هستند (فاتحی، ۱۳۸۷).

این پژوهش در طی مراحل زیر اجرا شده است:

## جمع‌آوری داده‌های خاک

در این پژوهش محل ۳۲۰ نقطه مشاهداتی به روش تصادفی انتخاب گردید و نمونه‌های خاک سطحی (۳۰-۰ سانتیمتر) از این نقاط در سطح ۲۴ هزار هکتار منطقه برداشت شد و درصد کربن آلی آن‌ها به عنوان متغیر هدف به روش سوزانیدن تر اندازه‌گیری گردید.

## متغیرهای محیطی کمکی

متغیرهای کمکی محیطی در این پژوهش به عنوان داده‌های ورودی پیش‌بینی کننده، با استفاده از مدل‌های رقومی ارتفاع با اندازه پیکسل ۱۰ متر، برخی مشتقات مراتب اول و دوم آن شامل ارتفاع، درصد شیب، تحدب، شاخص همواری دره با درجه

<sup>1</sup> model averaging (ensemble model)

<sup>2</sup> Equal weights averaging (EW)

<sup>3</sup> Bates-Granger or variance weighted averaging (VW)

<sup>4</sup> Granger-Rama

**There are no sources in the current document.**nathan averaging (GRA)

<sup>5</sup> Bayesian model averaging (BMA)

تفکیک بالا<sup>۶</sup>، شاخص رسوب<sup>۷</sup>، شاخص خیزی<sup>۸</sup>، عمق دره<sup>۹</sup>، فاصله عمودی تا شبکه آبراهه<sup>۱۰</sup> و شاخص پوشش گیاهی نرمال شده با استفاده از نرم افزار ایلویس نسخه ۳/۸ و ساگا نسخه ۲/۲ تهیه گردید.

در این مطالعه از مدل های خطی چندمتغیره، کوبیست<sup>۱۱</sup> و جنگل تصادفی<sup>۱۲</sup> برای پیش بینی مقدار کربن آلی خاک استفاده شد و در مرحله بعد با استفاده از مدل خطی تعمیم یافته<sup>۱۳</sup> مدل های مذکور ترکیب و یک مدل میانگین یا مدل تجمیعی به دست آمد و کلیه مدل ها با روش اعتبارسنجی دوجانبه مورد ارزیابی قرار گرفت.

## مدل های خطی چند متغیره

مدل رگرسیون خطی چند متغیره رابطه خطی بین متغیر هدف و حداقل دو متغیر پیش بینی کننده را نشان می دهد و اجرای آن نیازمند فرضیات مانند نرمال بودن توزیع و ثابت بودن واریانس متغیر هدف می باشد (McCullagh and Nelder, 1989).

## مدل کوبیست

مدل کوبیست به دلیل توانایی آن در کشف روابط غیرخطی داده ها دارای ساختار عامه پسندی است و مسئله پیش بینی های محدودی که در دیگر مدل های درختان طبقه بندی و رگرسیون اتفاق می افتد را ندارد. در مدل کوبیست، ابتدا داده ها را با توجه متغیر هدف و متغیرهای کمکی به زیر مجموعه هایی با ویژگی های مشابه، تفکیک می شوند. یک سری از قواعد زیرمجموعه های تفکیکی را معین می کنند. این قواعد به شکل سلسله مراتبی هستند. هر قاعده به شکل زیر است:

{ اگر شرایط واقعی است }

{ پس اجرای رگرسیون }

در غیر این صورت { اجرای قاعده بعدی }

شرایط ممکن است ساده و بر اساس یک متغیر کمکی باشد ولی اغلب تعدادی متغیر کمکی را شامل می شود. اگر شرایط واقعی باشد در مرحله بعد پیش بینی ویژگی خاک مورد نظر به وسیله اجرای مدل رگرسیون در درون هر زیر مجموعه تفکیکی انجام می گیرد. اگر شرایط درست (واقعی) نباشد بنابراین قاعده ای در گره بعدی درخت تعریف شده و توالی اگر، سپس، در غیر این صورت، تکرار می گردد. خروجی ها، حاصل از معادلات رگرسیونی هستند. مدل کوبیست براساس الگوریتم M5 کوینلان (۱۹۹۲) عمل می کند و در بسته نرم افزاری Cubist در محیط R اجرا گردید.

## مدل جنگل تصادفی

بریمن (۲۰۰۱) الگوریتم جنگل تصادفی را پیشنهاد نمود. در درختان طبقه بندی و رگرسیون معمولی، هر گره با استفاده از بهترین تفکیک کننده از میان همه ی متغیرهای کمکی منشعب می شود اما در جنگل تصادفی هر گره به طور تصادفی با استفاده از بهترین زیرمجموعه از بین زیرمجموعه های متغیرهای کمکی موجود در گره، منشعب می گردد. این یک راهبرد تا اندازه ای هوشمندانه است که در مقابل بیش برآزش، قوی عمل می کند و به آسانی قابل استفاده است. مدل جنگل تصادفی دارای دو پارامتر اصلی است، یکی تعداد متغیرها کمکی موجود در زیرمجموعه تصادفی هر گره (mtry) و دیگری تعداد درختان در جنگل (ntree) است که معمولاً به مقدار آن ها حساس نیست. از شاخص های آنتروپی مانند شاخص جینی (gini index) برای کمی سازی ناخالصی ها در هر گره استفاده می شود. در این روش ۳۷ درصد از مجموع کل داده ها در مرحله واسنجی مدل مورد استفاده قرار نمی گیرد که آنها را نمونه های بیرون از سبد (out of bag) می نامند. از این داده ها برای برآورد خطای پیش بینی مدل و به عبارت بهتر، اعتبارسنجی آن استفاده می شود که به آن خطای خارج از سبد (out of bag error) یا (OOB) می گویند.

<sup>6</sup> MRVBF

<sup>7</sup> Sediment index

<sup>8</sup> Topographic wetness index

<sup>9</sup> Valley Depth

<sup>10</sup> Vertical Distance to Channel Network

<sup>11</sup> cubist

<sup>12</sup> Random forest

<sup>13</sup> Generalized linear model

به همین دلیل معمولاً برای اعتبارسنجی مدل جنگل تصادفی از داده‌های مستقل استفاده نمی‌کنند. در این پژوهش برای اجرای مدل جنگل تصادفی از بسته‌ی نرم‌افزاری random forest در محیط R استفاده شد.

### میانگین مدل یا مدل تجمیعی

معمولاً در فرایند مدلسازی، پژوهشگران به منظور انتخاب بهترین مدل<sup>۱۴</sup>، اقدام به استفاده از مدل‌های مختلف و ارزیابی و مقایسه آن‌ها می‌کنند. بهترین مدل، یک مفهوم ذهنی است که وابسته به هدف پژوهشگر از مطالعه می‌باشد. یک روش جایگزین برای مفهوم "بهترین مدل" استفاده از روش‌های مدل میانگین (مدل تجمیعی) است (Padarian et al., 2014). مدل میانگین (مدل تجمیعی) ترکیب نمودن پیش‌بینی مدل‌های مختلف می‌باشد. اصول بنیادین این روش‌ها می‌تواند توسط مدل ساده زیر توضیح داده شود:

$$Y_i = \sum_{k=1}^K W_k X_{ik} \quad (2)$$

در این رابطه  $Y_i$  ترکیب خروجی مدل‌ها در نقطه  $i$ ، تعداد مدل‌های سهمیم در پیش‌بینی،  $X_{ik}$  خروجی  $k$  امین مدل و  $W_k$  ضریب وزنی اختصاص یافته به  $k$  امین مدل می‌باشد. در بیشتر روش‌های مدل میانگین، جمع اوزان برابر با یک است. روش‌های مدل میانگین مختلفی از ساده‌ترین حالت که اوزان مدل‌های مختلف یکسان در نظر گرفته می‌شود تا مدل‌های پیچیده‌ای مانند مدل میانگین بیزی می‌توان استفاده نمود. در این پژوهش از مدل‌های خطی تعمیم یافته<sup>۱۵</sup> برای ترکیب خطی مدل‌های مختلف و تولید مدل میانگین استفاده می‌شود که اصطلاحاً به این شیوه‌ی مدل میانگین، انبوهش مدل‌ها (stacking models) گفته می‌شود. چنین مدلی با استفاده از نرم افزار caret ensemble اجرا می‌گردد.

### اعتبارسنجی

در این پژوهش، از روش اعتبارسنجی دوجانبه<sup>۱۶</sup> استفاده شد. بدین منظور ۳۲۰ داده جمع‌آوری شده به دو دسته داده آموزشی (۲۴۰ نمونه) و اعتبارسنجی (۸۰ نمونه) تقسیم شد. خطای مدل‌ها بر اساس مرسوم‌ترین آن‌ها یعنی ریشه میانگین مربعات خطای تخمین<sup>۱۷</sup> محاسبه گردید:

$$RMSE = \sqrt{\frac{\sum_{i=1}^n [Z_{x_i}^* - Z_{x_i}]^2}{n}} \quad (3)$$

در این معادله،  $RMSE$  = ریشه میانگین مربعات خطای تخمین،  $Z_{x_i}^*$  = متغیر پیش‌بینی شده،  $Z_{x_i}$  = متغیر اندازه‌گیری شده و  $n$  = تعداد نمونه می‌باشد. یک مدل خوب، دارای حداقل مقدار ریشه میانگین مربعات خطای تخمین می‌باشد (بروس و همکاران، ۲۰۱۱). معیار دیگر برای تخمین دقت مدل، ضریب تعیین ( $R^2$ ) است که کارایی مدل در پیش‌بینی متغیر هدف با استفاده از متغیرهای کمکی را نشان می‌دهد.

$$R^2 = \frac{\sum_{i=1}^n [Z^*(x_i) - Z^*(\bar{x})]^2}{\sum_{i=1}^n [Z(x_i) - Z(\bar{x})]^2} \quad (4)$$

در معادله‌ی فوق  $Z(\bar{x})$  و  $Z^*(\bar{x})$  به ترتیب نشان دهنده‌ی میانگین داده‌های مشاهده‌ای و پیش‌بینی هستند.

<sup>14</sup> The best model

<sup>15</sup> Generalized linear models

<sup>16</sup> cross-validation

<sup>17</sup>Root Mean Squared Error

نتایج و بحث

در این تحقیق از ۹ متغیر کمکی شامل ارتفاع، درصد شیب، تحدب، شاخص همواری دره با درجه تفکیک بالا، شاخص رسوب، شاخص خیزی توپوگرافیک، عمق دره، فاصله عمودی تا شبکه آبراهه و شاخص پوشش گیاهی نرمال شده به عنوان متغیرهای محیطی کمکی و کربن آلی خاک به عنوان متغیر هدف استفاده شد. تمام مدل‌ها در ده تکرار اجرا شدند و ارزیابی درونی آن‌ها بر اساس اعتبار سنجی دو جانبه صورت گرفت و مدل بهینه بر اساس آن انتخاب گردید. همانطور که در جدول ۱ مشاهده می‌شود متغیرهای ارتفاع و شاخص پوشش گیاهی نرمال شده موثرترین متغیرهای ورودی در هر سه مدل هستند. ژائو و شی (۲۰۱۰) نشان دادند رابطه معناداری بین کربن آلی خاک و ارتفاع و شاخص پوشش گیاهی در مدل خطی چند متغیره وجود دارد که با نتایج پژوهش حاضر همخوانی دارد.

جدول ۱ - اهمیت نسبی متغیرهای کمکی در هر مدل

متغیرهای کمکی	رگرسیون خطی چند متغیره	مدل کویبست	جنگل تصادفی
ارتفاع	۱۰۰	۱۰۰	۲۵/۶۱
شاخص پوشش گیاهی نرمال شده	۵۴/۷۱۸	۳۲/۸۸	۱۴/۱۴
شاخص رسوب	۴۳/۲۲	۲۱/۲۳	۶/۶
فاصله عمودی تا شبکه آبراهه	۳۸/۷۰۱	۲۵/۳۴	۱۱
شاخص همواری دره با درجه تفکیک بالا	۳۶/۸۵۸	۶/۸۵	۱۲/۵۹
شیب	۲۵/۰۶۸	۱۴/۳۸	۷/۷۴
شاخص خیزی توپوگرافیک	۲۲/۳۰۱	۴۰/۴۱	۹/۵
تحدب	۲۲/۶۶۳	.	۷/۷۱
عمق دره	.	.	۷/۹۴

از میان مدل‌های منفرد، مدل جنگل تصادفی با ضریب تبیین (۰/۱۴) و ریشه میانگین مربعات خطا (۰/۵۷۵) دارای بیشترین صحت و دقت برای پیش بینی کربن آلی خاک است (جدول ۲). مدل بهینه جنگل تصادفی برای پیش بینی مقدار کربن آلی خاک شامل  $mtry = 3$  و  $tree = 500$  بود. اما نتایج اعتبارسنجی نشان می‌دهد مدل میانگین از نوع ترکیب خطی تعمیم یافته بیش از ۱۵ درصد تغییرات کربن آلی خاک در محدوده مورد مطالعه را توضیح می‌دهد که این مقدار در مدل‌های منفرد دیگر پایین تر است (ستون ضریب تبیین ( $R^2$ ) در جدول ۲ مشاهده شود). ریشه میانگین مربعات خطا در مدل میانگین از نوع ترکیب خطی تعمیم یافته (۰/۵۵۸) پایین تر از مدل‌های منفرد رگرسیون خطی چندمتغیره (۰/۵۸۶)، کویبست (۰/۵۸۱) و جنگل تصادفی (۰/۵۷۵) است (جدول ۲) که این نتایج خود بیانگر کارایی بیشتر مدل میانگین از نوع ترکیب خطی تعمیم یافته در پیش بینی مقدار کربن آلی خاک نسبت به سایر مدل‌های منفرد در این پژوهش می‌باشد. پادریان و همکاران (۲۰۱۴) از سه مدل کویبست، برنامه ریزی ژنتیکی و ماشین بردار پشتیبان و مدل‌های میانگین از نوع اوزان معادل و گرانگر-رامانتان برای پیش بینی ظرفیت رطوبت قابل استفاده خاک استفاده نمودند. نتایج آنها نشان داد که مدل‌های میانگین عملکرد بهتری از مدل‌های منفرد دارند

جدول ۲ - نتایج اعتبار سنجی مدل‌های مختلف برای پیش بینی کربن آلی خاک

نوع مدل مورد استفاده	R <sup>2</sup>	RMSE
خطی چند متغیره	۰/۰۹۶	۰/۵۸۶
کویبست	۰/۱۱۹	۰/۵۸۱
جنگل تصادفی	۰/۱۴۱	۰/۵۷۵
مدل میانگین از نوع ترکیب خطی تعمیم یافته	۰/۱۵۴	۰/۵۵۸



به طور کلی می‌توان بیان نمود استفاده از مدل‌های میانگین در نقشه‌برداری رقومی خاک باعث افزایش صحت و دقت نقشه‌های پیش‌بینی نسبت به مدل‌های منفرد حتی بهترین آن‌ها می‌شود.

## منابع

- Anderson J.A. (1995). An Introduction to Neural Networks. Cambridge, USA: The MIT Press.
- Breiman L., Friedman J.H, Olshen R.A., Stone & C.J.1984. Classification and regression trees.Wadsworth & Brooks. Monterey, CA.
- Brus D. J., Kempen B., and Heuvelink G. B. M. 2011. Sampling for validation of digital soil maps. Eur. J. Soil Sci., 62: 394-407.
- Cortes C. and Vapnik V.1995. Support-vector networks. Machine learning, 20(3): 273-297.
- Diks C.G.H., Vrugt J.A., 2010. Comparison of point forecast accuracy of model averaging methods in hydrologic applications. Stoch. Env. Res. Risk A., 24 (6): 809-820.
- Fatehi Sh. 2010. Semi-Detailed Soil Survey of Merak Plain in Kharkeh River Basin. Kermanshah Agricultural and Natural Resources Research and Education Center, Agricultural Research, Education and Extension Organization (AREEO), Kermanshah, Iran, 68p. (In Persian).
- Koza, J. (1992). Genetic Programming: On the Programming of Computers by Means of Natural Selection.The MIT Press.
- Malone B., McBratney A., Odgers N.P., & Minasny B.2014. Using model averaging to combine soil property rasters from legacy soil maps and from point data. Geoderma, 232-234: 34-44.
- McCullagh P. & Nelder J. 1989. Generalized Linear Model. Chapman & Hall, London.
- Padarian J., Minasny B., and McBratney A.B.2014. The evolving methodology for global soil mapping. In: Arrouays, D., McKenzie, N., Hempel, J., Richer de Forges, A., McBratney, A.B. (Eds.), GlobalSoilMap: Basis of the Global Spatial Soil Information System. CRC Press, pp. 215-220.
- Quinlan R.1992. Learning with continuous classes. Proceedings of the 5th Australian Joint Conference on Artificial Intelligence, Hobart, Tasmania, pp. 343-348.
- Zhao Y.C., and Shi X.Z.2010. Spatial Prediction and Uncertainty Assessment of Soil Organic Carbon in Hebei Province, China. In: Boettinger,J.L., Howell, D.W., More, A.C., Hartemink, A.E., Kienast-Brown, S. (Eds.), Digital Soil Mapping: Bridging Research, Environmental Application, and Operation. Springer, London, pp. 227-240.

### Using model averaging (ensemble model) in soil digital mapping

Sh. Fatehi

Research Assistant Professor, Soil and Water Research Department, Kermanshah Agricultural and Natural Resources Research and Education Center, AREEO, Kermanshah

#### Abstract

Using model averaging (ensemble model) is a strong way to improve the performance of digital soil properties modeling. but these techniques in digital soil mapping aren't regarded as a good alternative. This study aimed to evaluate and compare the prediction accuracy of three models include multivariate linear regression, Cubist and Random forest with a model averaging approach for mapping the soil organic carbon contents. First, in Merck sub-watershed in Kermanshah province with an area of 24,000 hectares, 320 soil samples were collected randomly and then soil organic carbon content was measured. Soil data were split in two sets, training (240 samples) and validation (80 sample). NDVI and the derived surface parameters of DEM were used as covariates in the fitted models. The root mean square error (RMSE) in model averaging (0.558) was lower than multivariable linear regression (0.586), Cubist (0.581) and Random Forest (0.575). In general, it can be said application of model averaging to digital soil mapping increased the accuracy of the prediction maps.

**Keywords:** digital soil modeling, data manning, ensemble model