

# مقایسه روش‌های رگرسیونی کلاسیک با شبکه عصبی و ماشین بردار پشتیبان در رده‌بندی منابع آب‌های زیرزمینی

اکرم حیدری<sup>۱</sup>، مهرداد نیپرست<sup>۲</sup>

تاریخ دریافت: ۹۸/۱/۲۷

تاریخ پذیرش: ۹۹/۳/۸

چکیده:

در عصر حاضر رده‌بندی داده‌ها به منظور تشخیص و پیش‌بینی وقایع، یکی از موضوعات بسیار مهم در علوم مختلف است. در دیدگاه سنتی علم آمار این رده‌بندی‌ها بر اساس روش‌های کلاسیک و بر پایه مدل‌های آماری از جمله رگرسیون لوژستیک امکان‌پذیر بود. در عصر حاضر که به عبارتی عصر انفجار اطلاعات نامیده می‌شود، در اکثر موارد با داده‌هایی مواجه هستیم که نمی‌توان توزیع دقیقی را برای آن‌ها یافت؛ از این رو استفاده از روش‌های داده کاوی و یادگیری ماشین که به مدل‌های از پیش تعیین شده نیاز ندارند، می‌تواند مثرتر باشد. در بسیاری از کشورها تشخیص دقیق نوع منابع آب‌های زیرزمینی، یکی از مسائل قابل توجه در زمینه علوم آب است. در این مقاله به مقایسه نتایج حاصل از رده‌بندی یک مجموعه داده مربوط به منابع آب‌های زیرزمینی با استفاده از روش‌های رگرسیونی، شبکه عصبی و ماشین بردار پشتیبان پرداخته‌ایم. نتایج این رده‌بندی‌ها نشان داد که روش‌های یادگیری ماشین در تشخیص دقیق نوع چشمه‌ها مؤثرتر بوده است.

واژه‌های کلیدی: شبکه عصبی<sup>۳</sup>، ماشین بردار پشتیبان<sup>۴</sup>، رگرسیون لوژستیک

## ۱ مقدمه

## ۲ رگرسیون دنباله‌رو تصویر ساز

برای بیان شبکه‌های عصبی لازم است ابتدا با مبانی تحت عنوان "رگرسیون دنباله‌رو تصویر ساز"<sup>۶</sup> ( $PPR$ ) آشنا شویم. به این منظور مطابق با شکل ۲ فرض کنید یک بردار ورودی  $X$  با  $p$  مؤلفه، یک بردار هدف  $y$  و نیز  $\omega_m$ ،  $m = 1, \dots, M$  واحد  $p$  مؤلفه‌ای از پارامترهای نامعلوم باشند. در این صورت مدل  $PPR$  به صورت زیر است؛

$$f(X) = \sum_{m=1}^M g_m(\omega_m^T X) \quad (1)$$

در این مدل جمعی، به جای ورودی‌ها، ویژگی‌های استخراج شده  $V_m = \omega_m^T X$  را قرار می‌دهیم و توابع  $g_m$  نیز نامشخص هستند و توسط برخی از روش‌های هموارسازی برآورد می‌شوند. متغیر عددی  $V_m = \omega_m^T X$  تصویر  $X$  را بر بردار واحد  $\omega_m$  تصویر می‌کند و تابع  $g_m(\omega_m^T X)$  یک تابع ستیغی<sup>۷</sup> در  $\mathbb{R}^p$  نامیده می‌شود. این تابع تنها در مسیر  $\omega_m$  تعریف شده است به همین

رگرسیون لوژستیک به عنوان یک ابزار رده‌بندی برای حالت‌هایی که متغیر پاسخ دو سطحی است، استفاده می‌گردد. تعمیم رگرسیون لوژستیک به حالت چند سطحی توسط گارلند و همکاران [۳]، منتل [۶] و تیل [۱۱] صورت گرفته است.

شبکه‌های عصبی و ماشین بردار پشتیبان نیز ابزارهای رده‌بندی برای مواردی هستند که توزیع دقیق داده‌ها مشخص نبوده و با یک مسئله ناپارامتری مواجه هستیم. الگوریتم شبکه عصبی نخستین بار توسط روزنبلات [۹] به صورت یک شبکه عصبی تک لایه با یک لایه پنهان که پرسپترون نامیده شد، مطرح گردید و الگوریتم  $SVM$  نیز نخستین بار توسط واپنیک و چروننکیس<sup>۵</sup> بیان شد. [۲]

<sup>۱</sup> کارشناسی ارشد آمار ریاضی، دانشگاه رازی، دانشکده علوم، گروه آمار

<sup>۲</sup> گروه آمار، دانشگاه رازی، دانشکده علوم، کرمانشاه، ایران

<sup>۳</sup>Neural Networks

<sup>۴</sup>Support Vector Machine

<sup>۵</sup>Vapnik and Chervonenkis

<sup>۶</sup>Projection Pursuit Regression

<sup>۷</sup>Ridge function

دلیل این مدل را مدل رگرسیون دنباله‌رو تصویرساز می‌نامند. برای ارزیابی مدل برازش داده‌شده، با توجه به این که شیوه آموزش در این روش یادگیری نظارتی<sup>۸</sup> است، با در نظر گرفتن داده آموزشی  $i = 1, \dots, N, (x_i, y_i)$  بر روی توابع  $g_m$  و بردارهای جهت  $\omega_m$ ، تابع خطای زیر را در نظر گرفته و آن را برآورد می‌کنیم؛

$$\sum_{i=1}^N [y_i - \sum_{m=1}^M g_m(\omega_m^T X)] \quad (2)$$

برای برآورد معادله فوق فرض می‌کنیم  $M = 1$  است. در این صورت باید آن را بر اساس  $\omega$ ، کمینه کنیم.

برای انجام این کار یک روش گاوس-نیوتون<sup>۹</sup> ابداع شد که روشی شبه-نیوتون<sup>۱۰</sup> است و هر بخش از مسیر شامل مشتقات دوم  $g$  است. در این صورت اگر  $\omega_{old}$  برآورد حاضر برای  $\omega$  باشد داریم؛

$$g(\omega^T x_i) \approx g(\omega_{old}^T x_i) + g'(\omega_{old}^T x_i)(\omega - \omega_{old})^T x_i \quad (3)$$

در این صورت می‌توان مقدار خطا را به صورت زیر بر اساس  $\omega$  برآورد کرد؛

$$\sum_{i=1}^N [y_i - g(\omega^T x_i)]^2 \approx \sum_{i=1}^N g'(\omega_{old}^T x_i)^2 \left[ (\omega_{old}^T x_i + \frac{y_i - g(\omega_{old}^T x_i)}{g'(\omega_{old}^T x_i)} - \omega^T x_i)^2 \right] \quad (4)$$

این عمل بردار ضرایب به‌روزرسانی شده  $\omega_{new}$  را به وجود می‌آورد و برآورد  $g$  و  $\omega$  تا زمانی که مسئله همگرا شود ادامه می‌یابد و در صورتی که در مدل  $PPR$ ،  $M > 1$  باشد، مدل در یک رویکرد روبه‌جلو ساخته می‌شود که در هر مرحله جفت  $(\omega_m, g_m)$  اضافه می‌گردد. در واقع می‌توان گفت که تعداد شروط  $M$  به عنوان بخشی از استراتژی روبه‌جلو به کار می‌رود و ساخت مدل زمانی متوقف می‌گردد که شرط بعدی موجب بهبود برازش مدل نگردد. به عبارت دیگر، استراتژی روبه‌جلو زمانی متوقف می‌شود که اجرای شرط بعدی که همان جفت  $(\omega_m, g_m)$  است، نتواند مدلی بهتر از مدل به‌دست آمده کنونی را به داده‌های موجود برازش دهد. [۵]

### ۳ شبکه عصبی

در این بخش نحوه عملکرد شبکه عصبی وانیلا<sup>۱۱</sup> که شبکه پرسرو با یک لایه پنهان و یا پرسپترون تک لایه نامیده می‌شود و نیز پرسپترون چند لایه<sup>۱۲</sup> ( $MLP$ ) را مورد بحث قرار می‌دهیم. شبکه‌های عصبی در واقع مدل‌های

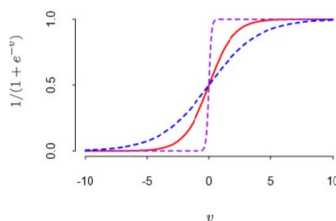
$$Z_m = \sigma(\alpha_m + \omega_m^T X), m = 1, \dots, M$$

$$T_k = \beta_{\cdot k} + \beta_k^T Z, k = 1, \dots, K \quad (5)$$

$$f_k(X) = g_k(T), k = 1, \dots, K$$

که در آن  $T = (T_1, \dots, T_K)$  و  $Z = (Z_1, \dots, Z_M)$  توجه کنید که در اینجا ضرایب  $\alpha$  مشابه با ضرایب  $\omega$  در مدل  $PPR$  است.

تابع فعال‌سازی  $\sigma(v)$  اغلب اوقات سیگموئید انتخاب می‌شود که به صورت  $\sigma(v) = \frac{1}{1 + \exp^{-v}}$  است. البته گاهی تابع پایه‌ای شعاعی گاوسی را برای  $\sigma(v)$  در نظر می‌گیرند که در این صورت آن را شبکه پایه‌ای شعاعی می‌نامند. اما در لایه‌های پنهان شبکه‌های عصبی که مقادیر  $Z_m$  در آن محاسبه می‌شود، اغلب از توابع سیگموئید استفاده می‌کنند. شکل ۱ نمایشی از تابع سیگموئید  $\sigma(sv)$  را نشان می‌دهد که در آن پارامتر  $s$ ، نرخ فعال‌سازی را کنترل می‌کند.  $s = \frac{1}{2}$  نمودار آبی رنگ و  $s = 10$  نمودار بنفش رنگ است. ملاحظه می‌کنید که  $s$  بزرگ در  $v = 0$  فعال‌سازی آن واحد (نورون-گره) را سخت می‌کند و  $\sigma(s(v - v_0))$  آستانه فعال‌سازی را از صفر به  $v_0$  انتقال می‌دهد. [۵]



شکل ۱. نمودار تابع سیگموئید (نمودار وسط)

شکل ۲ شمایی از یک شبکه عصبی پیشرو با یک لایه پنهان را نشان می‌دهد. یک ویژگی ورودی به نام اندازه اریبی<sup>۱۳</sup> وجود دارد که به هر یک از

<sup>8</sup>Supervised learning

<sup>9</sup>Gauss-Newton

<sup>10</sup>Quasi-Newton

<sup>11</sup>Vanilla

<sup>12</sup>Multi Layer Perceptron

<sup>13</sup>Bias

## ۴ برازش دادن شبکه‌های عصبی

شبکه‌های عصبی این امکان را برای کامپیوترها فراهم می‌آورند که به اطلاعات در مقیاس بزرگ دست یابند و نیز برای ساخت مدل‌های بسیار پیچیده به کار می‌روند. در واقع می‌توان گفت شبکه‌های عصبی «نبض» اصلی الگوریتم‌های یادگیری ماشین است. [۷] شبکه عصبی پارامترهای نامعلومی به نام وزن دارد که باید مقادیری را برای آن‌ها در نظر گرفت که برازش مدل به داده‌های آموزشی را به خوبی انجام دهد. می‌توان یک مجموعه کامل از وزن‌ها را با  $\theta$  نمایش داد که شامل مقادیر زیر است؛

$$\{\alpha_{\circ m}, \alpha_m; m = 1, \dots, M\}, \{\beta_{\circ k}, \beta_k; k = 1, \dots, K\} \quad (۷)$$

در این صورت برای رده‌بندی داده‌ها از یک خطای درجه دوم به نام آنتروپی-متقابل<sup>۱۵</sup> استفاده می‌کنیم؛ [۵]

$$R(\theta) = - \sum_{i=1}^N \sum_{k=1}^K y_{ik} \log f_k(x_i) \quad (۸)$$

در این صورت جداکننده متعلق به آن برابر است با؛

$$G(x) = \arg \max_k f_k(x)$$

با در نظر گرفتن تابع فعال‌سازی و تابع خطای آنتروپی-متقابل، مدل شبکه عصبی در واقع یک مدل رگرسیون لوژستیک خطی در واحدهای پنهان است که تمام پارامترهای آن با روش درستمایی ماکسیمم برآورد می‌شوند. حال می‌توان گفت که یک روش عمومی برای مینیمم کردن  $R(\theta)$  از طریق کاهش گرادیان<sup>۱۶</sup> است که تنظیمات انتشار پسر<sup>۱۷</sup> نامیده می‌شود. فرض کنید  $\sigma(\alpha_{\circ m} + \alpha_m^T x_i)$  و نیز با فرض (۵) از معادله (۵) داریم؛ [۵]

$$R(\theta) = \sum_{i=1}^N R_i = \sum_{i=1}^N \sum_{k=1}^K (y_{ik} - f_k(x_i))^2 \quad (۹)$$

$$\frac{\partial R_i}{\partial \beta_{km}} = -2(y_{ik} - f_k(x_i)) g'_k(\beta_k^T z_i) z_{mi},$$

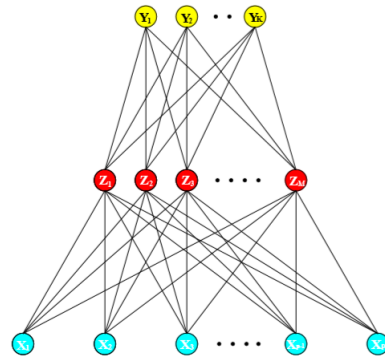
$$\frac{\partial R_i}{\partial \alpha_{ml}} = - \sum_{k=1}^K 2(y_{ik} - f_k(x_i)) g'_k(\beta_k^T z_i) \beta_{km} \sigma'(\alpha_m^T x_i) x_{il}. \quad (۱۰)$$

با استفاده از این مشتقات، یک بروزسانی کاهش گرادیان در  $(r+1)$  امین تکرار برابر است با؛

$$\beta_{km}^{(r+1)} = \beta_{km}^{(r)} - \gamma_r \sum_{i=1}^N \frac{\partial R_i}{\partial \beta_{km}^{(r)}},$$

$$\alpha_{ml}^{(r+1)} = \alpha_{ml}^{(r)} - \sum_{i=1}^N \frac{\partial R_i}{\partial \alpha_{ml}^{(r)}}. \quad (۱۱)$$

واحدهای موجود در لایه‌های خروجی و پنهان اضافه می‌شود و می‌توان آن را برابر با ثابت «۱» در نظر گرفت. این واحد آریبی در واقع عرض از مبدأهای  $\alpha_{\circ m}$  و  $\beta_{\circ m}$  موجود در معادله (۵) را می‌گیرد.



شکل ۲. شمایی از یک شبکه عصبی پیشرو با یک لایه پنهان

برای انجام محاسبات لازم برای  $Z_m$  ها به تابعی به نام توابع تشخیص<sup>۱۴</sup> نیاز است. این تابع در شبکه عصبی همان تابع سیگموئید است. مطابق با شکل ۱ و نیز معادله (۵) ملاحظه می‌کنید که نرخ فعال‌سازی سیگموئید به نرم  $\alpha_m$  بستگی دارد به گونه‌ای که اگر  $\|\alpha_m\|$  خیلی کوچک باشد، بخش خطی تابع فعال‌سازی را فراهم می‌آورد. با توجه به این مطالب می‌توان گفت که شبکه عصبی بسیار شبیه به مدل  $PPR$  است و تفاوت آن‌ها در این است که مدل  $PPR$  از توابع ناپارامتری  $g_m(v)$  استفاده می‌کند درحالی‌که شبکه عصبی از یک تابع بسیار ساده‌تر بر اساس  $\sigma(v)$  با سه پارامتر آزاد نرخ فعال‌سازی،  $\alpha$  و  $\beta$  استفاده می‌کند. در این صورت اگر از دیدگاه مدل‌های  $PPR$  به شبکه عصبی نگاه کنیم داریم؛

$$g_m(\omega_m^T X) = \beta_m \sigma(\alpha_{\circ m} + \alpha_m^T X)$$

$$= \beta_m \sigma(\alpha_{\circ m} + \|\alpha_m\| (\omega_m^T X)) \quad (۶)$$

طوری که  $\omega_m = \frac{\alpha_m}{\|\alpha_m\|}$ ،  $m$  امین بردار واحد است. در واقع مزیت شبکه عصبی بر مدل  $PPR$  را می‌توان این‌گونه بیان نمود:

با توجه به این که  $\sigma_{\beta, \alpha_{\circ}, s(v)} = \beta \sigma(\alpha_{\circ} + sv)$  نسبت به یک  $g(v)$  ناپارامتری، پیچیدگی کم‌تری دارد، لذا شبکه عصبی می‌تواند برای تحلیل مدل‌هایی با تعداد ورودی‌های بسیار بالا، حتی تا  $10^6$  متغیر نیز مورد استفاده قرار می‌گیرد. این در حالی است که مدل  $PPR$  از قیود کم‌تری استفاده می‌کند به‌عنوان مثال  $M = 5$  یا  $10$ . [۵]

<sup>14</sup>Identity functions

<sup>15</sup>Cross-Entropy

<sup>16</sup>Gradian decay

<sup>17</sup>Back-propagation

درون هر لایه، جمع وزن دهی شده به دست آمده از نورون‌های لایه قبلی را محاسبه کرده و در صورتی که مقدار آن از یک حد آستانه معین بیش تر شود، آن واحد (نورون) فعال خواهد شد. [۷]

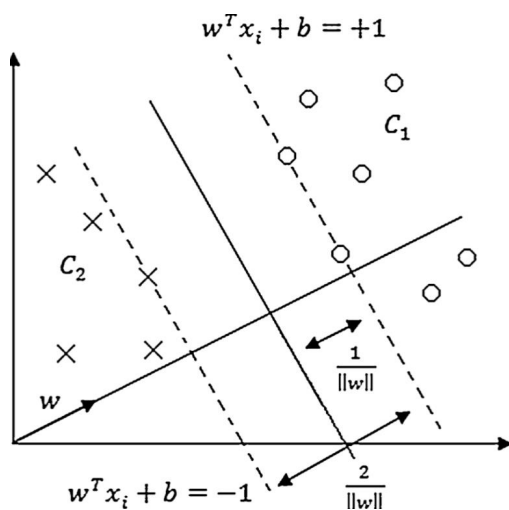
## ۵ ماشین بردار پشتیبان

یکی از شاخه‌های یادگیری ماشین، SVM است و برای رده‌بندی داده‌ها به کار می‌رود که بنا بر تعداد سطوح متغیر پاسخ به دو حالت دو رده‌ای و چند رده‌ای تقسیم می‌گردد. فرآیند جداسازی رده‌ها توسط توابعی به نام تابع تصمیم انجام می‌شود. اگر داده‌ها در فضای ورودی اولیه به صورت خطی جداپذیر باشند تابع تصمیم را خطی خوانده و SVM را سخت-حاشیه می‌نامند. تابع تصمیم در حالت دو رده‌ای را به صورت زیر نمایش می‌دهند

$$y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1, \quad i = 1, 2, \dots, N$$

به طوری که  $\mathbf{w}$  یک بردار  $p$  بعدی که  $p$  بیان‌گر بعد داده مشاهده شده  $\mathbf{x}$  (داده مشاهده شده  $i$  ام را با  $\mathbf{x}_i$  نمایش می‌دهیم) و  $b$  اندازه اریبی و  $N$  تعداد داده‌های آموزشی است. در این صورت تخصیص داده  $\mathbf{x}$  به هر یک از دو رده به این صورت انجام می‌شود؛

$$D(\mathbf{x}) \begin{cases} > 0 & \text{متعلق به رده یک} \\ < 0 & \text{متعلق به رده دو} \end{cases}$$



شکل ۳. SVM جداپذیر خطی

طوری که  $\gamma_r$  نرخ یادگیری نامیده می‌شود. می‌توان معادلات (۱۰) را به صورت زیر بازنویسی نمود؛

$$\frac{\partial R_i}{\partial \beta_{km}} = \delta_{ki} z_{mi}$$

$$\frac{\partial R_i}{\partial \alpha_{ml}} = s_{mi} x_{il} \quad (12)$$

طوری که  $\delta_{ki}$  و  $s_{mi}$  «خطاهای» مدل اصلاح شده به ترتیب در واحدهای خروجی و پنهان هستند و داریم؛

$$s_{mi} = \sigma'(\alpha_m^T x_i) \sum_{k=1}^K \beta_{km} \delta_{ki} \quad (13)$$

این معادلات را معادلات انتشار پس‌رو می‌نامند و با استفاده از به‌روزرسانی‌ها در معادلات (۱۱) یک الگوریتم دو گذرگاهی<sup>۱۸</sup> با دو گذرگاه پیشرو<sup>۱۹</sup> و پس‌رو<sup>۲۰</sup> ایجاد می‌گردد به نحوی که در گذرگاه پیشرو، وزن‌های اصلاح شده ثابت شده‌اند و مقادیر برآورد شده  $\hat{f}_k(x_i)$  از معادله (۵) به دست می‌آید و نیز در گذرگاه پس‌رو، خطاهای  $\delta_{ki}$  محاسبه می‌شوند. انتشار پس‌رو از طریق معادله (۱۳) خطاهای  $s_{mi}$  را به دست می‌آورد، سپس هر دو مجموعه خطا را برای محاسبه گرادیان‌ها برای به‌روزرسانی‌ها در معادلات (۱۱) از طریق معادلات (۱۲) مورد استفاده قرار می‌گیرند. این روش در واقع همان روش انتشار پس‌رو است که قانون دلنا نیز نامیده می‌شود. به‌روزرسانی‌ها در معادلات (۱۱) نوعی دیگر از یادگیری دسته‌ای<sup>۲۱</sup> است. اغلب نرخ یادگیری  $\gamma_r$  برای یادگیری دسته‌ای را یک مقدار ثابت در نظر می‌گیریم طوری که در هر مرحله، تابع خطا را مینیمم کند، به صورتی که هرگاه داشته باشیم  $r \rightarrow \infty$  آنگاه  $\gamma_r \rightarrow 0$ . بنابراین همگرایی در صورتی حتمی خواهد بود که

$$\gamma_r \rightarrow 0, \quad \sum_r r \gamma_r = \infty, \quad \sum_r \gamma_r^2 < \infty$$

(به‌عنوان مثال برای  $\gamma_r = \frac{1}{r}$  صدق می‌کند). [۵]

نکته‌ای که باید به آن توجه داشت این است که اگر وزن‌ها نزدیک به صفر باشند نمودار سیگموئید موجود در شکل ۱ تقریباً خطی است، در نتیجه شبکه عصبی به مدل خطی نزدیک می‌شود. معمولاً مقادیر شروع کننده وزن‌ها برای لایه اول، مقادیر تصادفی نزدیک به صفر انتخاب می‌شوند؛ به همین دلیل هم مدل با شکلی نزدیک به خطی شروع می‌شود و با افزایش وزن‌ها غیر خطی می‌شود. استفاده از وزن‌های صفر واقعی منجر به مشتقات صفر شده و الگوریتم هرگز حرکت نمی‌کند، [۵] زیرا علت اصلی فعال شدن شبکه در واقع فعال بودن اتصالات بین نورون‌های موجود در دو لایه مختلف است که این اتصالات، همان وزن‌ها هستند و توابع فعال‌سازی موجود در واحدهای

<sup>18</sup>Two-pass algorithm

<sup>19</sup>Forward-pass

<sup>20</sup>Backward-pass

<sup>21</sup>Batch learning

به طوری که

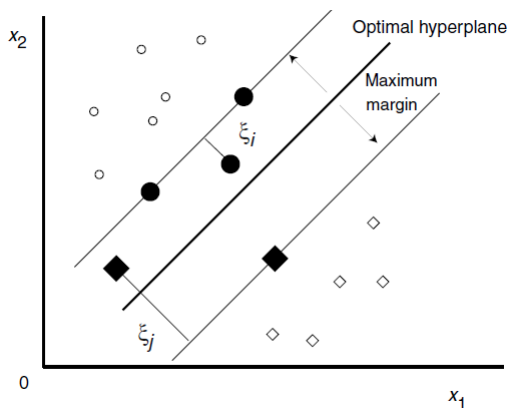
$$\sum_{i=1}^N \alpha_i y_i = 0, \quad \alpha_i \geq 0 \quad i = 1, 2, \dots, N$$

اما گاهی داده‌ها در فضای ورودی به صورت خطی قابل جداسازی نیستند. در این شرایط با دو حالت کلی مواجه هستیم  
**حالت اول**، در این حالت با افزودن متغیر ضعف  $\xi_i$  به تابع تصمیم، به صورت

$$y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i \quad , \quad i = 1, 2, \dots, N$$

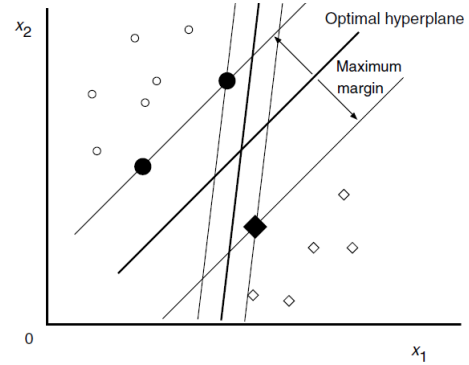
مسئله مینیم سازی (۱۵) را داریم و SVM را نرم-حاشیه می‌نامند. متغیر ضعف برای داده‌هایی که درون حاشیه بین دو ابرصفحه قرار می‌گیرند به کار می‌رود. زیرا برای سایر داده‌ها که به درستی رده‌بندی شده‌اند،  $\xi_i = 0$  است. با توجه به این که مرز تصمیم دارای معادله  $D(\mathbf{x}) = 0$  و ابرصفحه مثبت دارای معادله  $D(\mathbf{x}) = +1$  و نیز ابرصفحه منفی معادله  $D(\mathbf{x}) = -1$  را دارا است، لذا مطابق با شکل ۵ می‌توان برای داده آموزشی  $\mathbf{x}_i$  دو حالت زیر را در نظر گرفت  
۱- اگر  $1 < \xi_i < 1$  باشد، داده  $\mathbf{x}_i$ ، ماکسیم حاشیه را ندارد اما همچنان به درستی قابل رده‌بندی است،

۲- اگر  $\xi_i > 1$  باشد، (مانند  $\xi_j > 1$ )، در این صورت باید گفت داده ورودی به وسیله ابرصفحه بهینه به درستی رده‌بندی نشده است. [۱]



شکل ۵. استفاده از متغیر ضعف در رده‌بندی داده‌ها در حالت دو رده‌ای

$$\begin{aligned} \text{minimize } Q(\mathbf{w}, b, \xi, \alpha, \beta) = & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \xi_i \\ & - \sum_{i=1}^N \alpha_i (y_i (\mathbf{w}^T \mathbf{x}_i + b) - 1 + \xi_i) - \sum_{i=1}^N \beta_i \xi_i \end{aligned} \quad (15)$$



شکل ۴. ابرصفحه جداکننده بهینه در حالت دوبعدی

در شکل ۳ نمایی از جداپذیری خطی در حالت دو رده‌ای را مشاهده می‌کنیم طوری که خط چین‌ها ابرصفحه‌های مثبت و منفی هستند. [۱] مطابق با شکل ۴ انتخاب این تابع تصمیم برای جداسازی دو رده باید به گونه‌ای انجام شود که بیشینه حاشیه ممکن ( $h$ ) بین نزدیک‌ترین داده‌های مربوط به دو رده مختلف ایجاد گردد. اگر  $\mathbf{x}^+$  و  $\mathbf{x}^-$  به ترتیب داده‌هایی بر دو ابرصفحه مثبت و منفی باشند، حاشیه بین دو ابرصفحه برابر

$$h = |\mathbf{x}^+ - \mathbf{x}^-| = \frac{2}{\|\mathbf{w}\|}$$

خواهد بود. در این صورت مسئله بهینه‌سازی محدب زیر را داریم

$$\begin{aligned} \text{minimize } Q(\mathbf{w}, b, \alpha) = & \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_{i=1}^N \alpha_i (y_i (\mathbf{w}^T \mathbf{x}_i + b) - 1) \end{aligned} \quad (14)$$

به طوری که

$$(\alpha_i \geq 0), \quad \alpha = (\alpha_1, \alpha_2, \dots, \alpha_M)^T$$

که در آن  $\alpha_i$  ها ضرایب لاگرانژ نامفی مربوط به  $\mathbf{x}_i$  ها هستند. تحت شرایط "کاروش-کان-توکر"<sup>۲۲</sup> زیر داریم

$$\begin{aligned} \frac{\partial Q(\mathbf{w}, b, \alpha)}{\partial \mathbf{w}} = \mathbf{w} - \sum_i \alpha_i y_i \mathbf{x}_i = 0 & \Rightarrow \mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i \\ \frac{\partial Q(\mathbf{w}, b, \alpha)}{\partial \mathbf{b}} = - \sum_i \alpha_i y_i = 0 & \Rightarrow \sum_i \alpha_i y_i = 0 \end{aligned}$$

$\alpha_i$  هایی که مثبت دارند، بردارهای پشتیبان مسئله را تشکیل می‌دهند و روی ابرصفحه‌ها قرار دارند بنابراین با جایگذاری مقادیر فوق در معادله (۱۴) خواهیم داشت

$$\text{maximize } Q(\alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$$

<sup>22</sup>Karush-Kuhn-Tucker(KKT)

<sup>23</sup>Slack Variable

شکل ۶ نمایی از انتقال داده‌ها توسط کرنل چندجمله‌ای، از فضای دوبعدی به فضای مشخصه سه‌بعدی و جداسازی خطی داده‌ها در این فضا را نمایش می‌دهد.

این ایده به بیش از دو رده نیز قابل تعمیم است. برای جداسازی رده‌ها در این حالت می‌توان از چهار ماشین بردار پشتیبان متفاوت تحت عنوان یک در مقابل همه (OAA)، زوجی (PSVM)، اصلاح خطای کدهای خروجی (ECOC) و همه در یک‌زمان بهره گرفت. [۵]

## ۶ رگرسیون لوژستیک

با توجه به این که متغیر پاسخ کیفی و دارای چند سطح است لذا این مدل به منظور تخمین احتمال وقوع یک رویداد خاص (احتمال تعلق بردار داده  $\mathbf{x}$  به رده  $i$ ،  $i = 1, \dots, n$ )، به کار می‌رود، به طوری که  $n$  تعداد رده‌ها است. برای تخصیص داده  $\mathbf{x}$  به یکی از دو رده  $i$  و  $j$  از بین  $n$  رده موجود، از شاخصی تحت عنوان نسبت بخت‌ها (OR) به صورت زیر استفاده می‌شود

$$OR = \frac{\frac{\pi_i}{1-\pi_i}}{\frac{\pi_j}{1-\pi_j}}$$

به طوری که  $\pi_i = Pr(G = i | \mathbf{X} = \mathbf{x})$  و  $\pi_j = Pr(G = j | \mathbf{X} = \mathbf{x})$ ، در این صورت  $\pi_i$  را احتمال تعلق بردار داده  $\mathbf{x}$  به رده  $i$  می‌نامند. بر اساس مقدار حاصل شده  $OR$ ،  $\mathbf{x}$  را به یکی از دو رده  $i$  و  $j$  نسبت می‌دهند طوری که اگر  $OR > 1$  باشد آنگاه  $\mathbf{x}$  عضو رده  $i$ ، اگر  $OR < 1$  باشد آنگاه  $\mathbf{x}$  متعلق به رده  $j$  و اگر  $OR = 1$  باشد،  $\mathbf{x}$  غیرقابل رده‌بندی است. با در نظر گرفتن تابع ربط کانونی لجیت داریم

$$\log \frac{\pi_i}{\pi_n} = \frac{Pr(G = i | \mathbf{X} = \mathbf{x})}{Pr(G = n | \mathbf{X} = \mathbf{x})} = \beta_{i0} + \beta_i \mathbf{x}; i = 1, \dots, n-1$$

توجه شود که مقدار فوق در واقع همان آماره  $OR$  برای دو رده است و رده مخرج نیز قراردادی است. [۵]

## ۷ بحث و نتیجه‌گیری

در این پژوهش داده‌های مربوط به ۴۲ منبع آب زیرزمینی، به منظور بررسی و مقایسه رده‌بندی چشمه‌ها به دو دسته تحول یافته (A) و تحول

در این حالت نیز با استفاده از شرایط KKT داریم

$$\frac{\partial Q(\mathbf{w}, b, \xi, \alpha, \beta)}{\partial \mathbf{w}} = \mathbf{0} \Rightarrow \mathbf{w} - \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i = \mathbf{0} \Rightarrow \mathbf{w} = \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i$$

$$\frac{\partial Q(\mathbf{w}, b, \xi, \alpha, \beta)}{\partial b} = 0 \Rightarrow \sum_{i=1}^N \alpha_i y_i = 0$$

$$\frac{\partial Q(\mathbf{w}, b, \xi, \alpha, \beta)}{\partial \xi} = \mathbf{0} \Rightarrow C - \alpha_i - \beta_i = \mathbf{0} \Rightarrow \alpha_i + \beta_i = C$$

به طوری که  $C$  را پارامتر هزینه خوانده و رابطه بین ماکسیم سازی حاشیه بین دو رده و مینیم سازی خطای رده‌بندی را ارزیابی می‌کند. با جایگذاری مقادیر فوق در معادله (۱۵) داریم

$$\text{maximize } Q(\alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1, j \neq i}^N \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$$

$$\sum_{i=1}^N \alpha_i y_i = 0 \quad \text{و} \quad 0 \leq \alpha_i \leq C \quad i = (1, 2, \dots, N)$$

در این حالت نیز که تنها تفاوت آن با SVM سخت-حاشیه در کران بالای  $\alpha_i$  است، تابع تصمیم دقیقاً برابر با این تابع در حالت سخت-حاشیه است

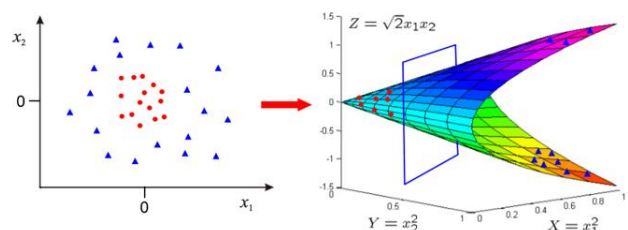
$$D(\mathbf{x}) = \sum_{i \in S} \alpha_i y_i \mathbf{x}_i^T \mathbf{x} + b$$

که در آن  $S$  مجموعه‌ای از  $\mathbf{x}_i$ ها با ضرایب لاگرانژ نامنفی است و برای رده‌بندی مشاهده  $\mathbf{x}$  داریم

$$D(\mathbf{x}) \begin{cases} > 0 & \text{متعلق به رده یک} \\ < 0 & \text{متعلق به رده دو} \end{cases}$$

**حالت دوم**، مربوط به زمانی است که حتی با اضافه نمودن متغیر ضعف به مسئله باز هم داده‌ها به صورت خطی قابل رده‌بندی نیستند؛ در این حالت می‌توان داده‌ها را توسط توابعی به نام کرنل که خود دارای انواع مختلفی است به فضایی با ابعاد  $l > p$  به نام فضای مشخصه انتقال داد طوری که داده‌ها در این فضای جدید، جداپذیر خواهند بود. این عمل را حیلۀ کرنل می‌نامند.

$$\Phi: \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \rightarrow \begin{pmatrix} x_1^2 \\ x_2^2 \\ \sqrt{2}x_1x_2 \end{pmatrix} \quad \mathbb{R}^2 \rightarrow \mathbb{R}^3$$



شکل ۶. نمایی از انتقال داده‌ها به فضای مشخصه

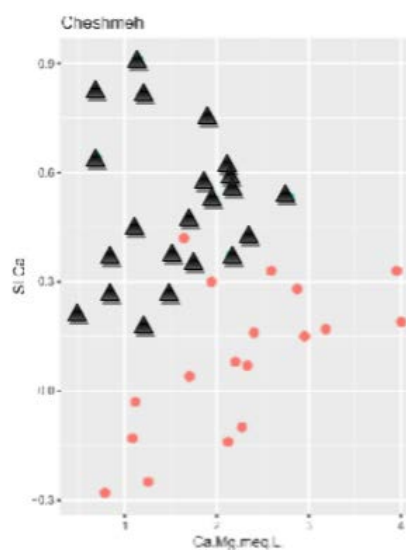
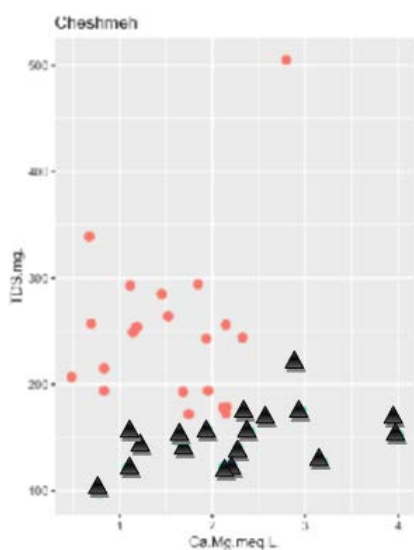
جدول ۱: رده‌بندی مجموعه داده منابع آب با استفاده از روش SVM

دقت	تعداد رده‌بندی اشتباه	SV	Coef <sup>o</sup>	degree	gamma	cost	هسته
۱	۰	۱۸	-	-	۰/۲۵	۱	radial
۰/۹۵	۲	۱۱	-	-	۰/۲۵	۱	linear
۰/۸۶	۶	۱۹	۰	۳	۰/۲۵	۱	polynomial
۰/۹۵	۲	۱۵	۰	-	۰/۲۵	۱	sigmoid

جدول ۲: ضرایب رگرسیونی برای دو مدل لجیت و پرویت

پرویت			لوجیت			متغیر مستقل
p-value	ErrorStd.	Estimate	p-value	ErrorStd.	Estimate	
< 2e-16	6.656e+07	5.562e+15	< 2e-16	6.656e+07	1.727e+15	Intercept
< 2e-16	1.382e+07	1.645e+15	< 2e-16	1.382e+07	5.991e+14	Ca/Mg
< 2e-16	5.636e+07	-3.785e+15	< 2e-16	5.636e+07	-1.248e+15	SI.Ca
< 2e-16	1.846e+05	-2.401e+13	< 2e-16	1.846e+05	-8.978e+12	TDS
< 2e-16	1.154e+05	-7.521e+12	< 2e-16	1.154e+05	-2.234e+12	EC
۱۵۴/۱۷۰			۸۲/۸۷			AIC
۰/۹۵۲۳۸۱			۰/۹۷۶۱۹۰۵			accuracy

نیافته (B) با استفاده از روش رگرسیون لوژستیک و دو روش ماشین بردار پشتیبان و شبکه عصبی به کاررفته است. در ادامه جدول‌ها و نتایج حاصل از این بررسی ارائه شده است.



شکل ۸. رده‌بندی منابع آب با در نظر گرفتن دو متغیر Ca/Mg و TDS

شکل ۷. رده‌بندی منابع آب با در نظر گرفتن دو متغیر Ca/Mg و SI.Ca

جدول ۳: دقت به‌دست‌آمده از سه روش شبکه عصبی، SVM و رگرسیون

روش	دقت
شبکه عصبی	۱
SVM	۱
رگرسیون (لوژستیک)	۰.۸۷

باید یک پارامتر جریمه را در قالب عمل ارزشیابی متقابل به مسئله اعمال کرد، که البته در جریان این پژوهش، این امر در نظر گرفته شده است. مطابق با نتایج حاصل شده از جدول ۳ می‌توان گفت با استفاده از روش‌های آموزش ماشین به راحتی و با میزان دقت ۱ می‌توان عمل رده‌بندی را انجام داد. پس از انجام محاسبات با استفاده از شبکه عصبی، ملاحظه می‌کنیم که تعداد لایه‌های پنهان برابر با ۱۳، حد آستانه برابر با ۱/۰ و تعداد تکرارهای لازم برای همگرا شدن مسئله برابر با ۳۵۰ مرحله است. مطابق با جدول ۱ دقت محاسبات با روش SVM با استفاده از چهار هسته نام برده در جدول بالا بوده طوری که بالاترین دقت مربوط به هسته رادیال با ۱۸ بردار پشتیبان و با دقت ۱ است. مطابق با جدول ۲ دقت در روش‌های رگرسیونی نیز بالا است. این میزان دقت در رگرسیون لوژستیک با تابع ربط لجیت که دارای ملاک آکائیکه  $۸۲/۰۸۷$  و دقت ۰.۸۷ است نسبت به تابع ربط پروبیت با  $AIC = ۱۵۴/۱۷$  و دقت ۰.۸۵، بالاتر است.

در شکل‌های ۷ و ۸ نقاط دایره شکل نشان‌دهنده داده‌های رده A با کد ۱ و نقاط مثلث شکل نشان‌دهنده داده‌های رده B با کد ۰ هستند. مطابق با این دو شکل ملاحظه می‌کنید، زمانی که الگوریتم شبکه عصبی برای رده‌بندی داده‌ها تنها دو متغیر  $Ca/Mg$  و  $SI.Ca$  در شکل ۷ و یا دو متغیر  $Ca/Mg$  و  $TDS$  در شکل ۸ را برای جداسازی دو رده A و B در نظر می‌گیرد، جداسازی آن‌ها با استفاده از روش شبکه عصبی بسیار دقیق صورت گرفته است البته این امر برای تمام جفت متغیرها در این مسئله برقرار است.

دو عامل مهم در تنظیم و افزایش یا کاهش نرخ خطا در شبکه‌های عصبی عبارت است از تعداد لایه‌های پنهان و تعداد واحدهای موجود در هر لایه، [۷] طوری که هر میزان تعداد لایه‌های پنهان بیشتر باشد، انعطاف‌پذیری شبکه افزایش یافته و دقت محاسبات افزایش می‌یابد، از طرفی این تعداد را نمی‌توان تا هر میزان دلخواه افزایش داد زیرا ممکن است مسئله به جواب درست همگرا نشود. بنابراین برای یافتن تعداد لایه‌های پنهان بهینه، مانند الگوریتم SVM،

## مراجع

- [1] Abe, S. (2010). *Support Vector Machines for Pattern Classification* (pp. 331-341). Springer, London.
- [2] Boser, B. E., Guyon, I. M., and Vapnik, V. N. (2003). A training algorithm for optimal margin classifiers. *In Proceedings of the 5th Annual ACM Workshop on Computational Learning Theory*, 144-152.
- [3] Gurland, J., Lee, I., and Dahm, P. A. (1960). Polychotomous quantal response in biological assay. *Biometrics*, **16**(3), 382-398.
- [4] Hastie, T., Tibshirani, R., and Friedman, J. (2009). *The elements of statistical learning: prediction, inference and data mining*. Springer-Verlag, New York.
- [5] Hastie, T., Tibshirani, R., Friedman, J., and Franklin, J. (2005). The elements of statistical learning: data mining, inference and prediction. *The Mathematical Intelligencer*, **27**(2), 83-85.
- [6] Mantel, N. (1966). Models for complex contingency tables and polychotomous dosage response curves. *Biometrics*, **22**(1), 83-95.

- [7] Muller, A. C., and Guido, S. (2017). *Introduction to machine learning with Python: a guide for data scientists*. O'Reilly Media.
- [8] Robbins, H., and Monro, S. (1951). A stochastic approximation method. *The annals of mathematical statistics*, **22**(3), 400-407.
- [9] Rosenblatt, F. (1958). The perceptron: a probabilistic model for information storage and organization in the brain. *Psychological review*, **65**(6), 386-408.
- [10] Rumelhart, D. E., Hinton, G. E., and Williams, R. J. (1986). Learning representations by back-propagating errors. *Nature*, **323**, 533-536.
- [11] Theil, H. (1969). A multinomial extension of the linear logit model. *International economic review*, **10**(3), 251-259.
- [12] Widrow, B., and Hoff, M. E. (1960). Adaptive switching circuits. 1960 IRE WESCON convention record (pp. 96-104). *New York: IRE. Reprinted in Anderson and Rosenfeld (1988)*, 126-134.