

# الگوریتم CLIP4

حسن شجاعی مند

دانشجوی کارشناسی ارشد نرم افزار

دانشکده مهندسی کامپیوتر

دانشگاه صنعتی شریف

shojaee@ce.sharif.edu

**خلاصه:** در این مقاله، می‌خواهیم الگوریتم CLIP4 را مورد بررسی قرار دهیم. این الگوریتم از نوع الگوریتم‌های یادگیری ماشین استقرایی می‌باشد، که برای انجام عمل دسته‌بندی<sup>۱</sup> مورد استفاده قرار می‌گیرد. الگوریتم ابتدا داده‌ها را به زیرمجموعه‌هایی با استفاده از یک ساختار درختی تقسیم‌بندی می‌کند و سپس قواعد تولید<sup>۲</sup> را تنها با استفاده از زیرمجموعه‌هایی که در گره‌های برگ قرار دارند، ایجاد می‌کند. ویژگی منحصر به فرد الگوریتم، امکان تولید قواعدی می‌باشد که از نامساوی استفاده می‌کنند. کارایی الگوریتم با استفاده از یکسری منبع داده محک‌زنی<sup>۳</sup>، با بقیه الگوریتم‌های یادگیری ماشین مقایسه خواهد شد.

## ۱- مقدمه

در این مقاله می‌خواهیم الگوریتم CLIP4 که در [1] مطرح شده است را مورد بررسی قرار دهیم. این الگوریتم یکی از الگوریتم‌های یادگیری ماشین به حساب می‌آید که از روش هدایت‌شده<sup>۴</sup>، استفاده می‌کند. این تنها الگوریتمی است که توانایی تولید قواعد نامساوی را دارا می‌باشد.

اولین الگوریتم CLIP<sup>۵</sup> در سال ۱۹۹۵ با عنوان CLIP2 در [2] ارائه شد. این الگوریتم بعداً در CLIP3 بهبود یافت [3]. الگوریتم CLIP4، تغییراتی در دو الگوریتم قبلی بوجود آورده است که باعث دقت، کارایی و انعطاف‌پذیری بیشتری می‌شود. این الگوریتم یک الگوریتم ترکیبی است که ایده‌های موجود در خانواده الگوریتم‌های یادگیری ماشین را با هم ترکیب کرده است: الگوریتم‌های مبتنی بر قاعده و الگوریتم‌های درخت تصمیم. تفاوت اصلی بین این الگوریتم و الگوریتم‌های موجود در این دو خانواده، استفاده زیاد آن از مسئله پوشش مجموعه<sup>۶</sup> می‌باشد. الگوریتم دیگری که به صورت ترکیبی از الگوریتم‌های مبتنی بر قاعده و الگوریتم‌های درخت تصمیم استفاده می‌کند، الگوریتم CN2 [5] می‌باشد که از یک معیار انتروپی شبیه به آنچه در الگوریتم‌های درخت تصمیم استفاده می‌شود، استفاده می‌کند. چند مورد در [1] به عنوان فاکتورهای مهمی که می‌توانند در یک سیستم یادگیری ماشین هدایت شده مطرح باشند، عنوان شده است:

<sup>1</sup> classification

<sup>2</sup> Production rule

<sup>3</sup> benchmarking

<sup>4</sup> supervised

<sup>5</sup> cover learning using integer linear programming

<sup>6</sup> set covering

- دسته‌بندی دقیق<sup>7</sup>: قواعد تولید شده، باید توانایی دسته‌بندی درست رکوردهایی که در آینده مطرح خواهند شد را داشته باشد و باید بتواند در حضور نویز در داده‌ها کار خود را بدرستی انجام دهد.
  - قواعد ساده: قواعد باید تا حد ممکن مختصر باشند. چرا که در اینصورت قابل فهم‌تر خواهند بود.
  - تولید کارای قواعد: الگوریتم باید توانایی کار با یک منبع داده بزرگ را نیز داشته باشد.
  - انعطاف‌پذیری: الگوریتم باید حوزه وسیعی از کاربردها را شامل شود. بتواند انواع داده‌های مختلف، تعداد فیلدهای زیاد و فیلدهای با داده‌های خطادار را تحت پوشش قرار دهد.
- در ادامه، ابتدا الگوریتم CLIP4 را با جزئیات و بر اساس شبه‌کد ارائه شده برای آن بررسی خواهیم کرد و سپس مثالی از کاربرد آنرا خواهیم دید و در ضمن مثال جزئیات کارکرد قسمتهای مختلف الگوریتم را توضیح خواهیم داد. سپس ویژگی‌های الگوریتم را بررسی کرده و در پایان نتایج مقایسه کارایی این الگوریتم با الگوریتم‌های موجود با استفاده از یکسری منبع داده محک‌زنی، ارائه خواهد شد.

## ۲- الگوریتم CLIP4

توصیف مفاهیم یادگیری، به صورت قواعد if-then یا درخت‌های تصمیم، از مطرح‌ترین روش‌های یادگیری ماشین می‌باشند. هر الگوریتم یادگیری ماشین هدایت‌شده قواعد دسته‌بندی را با استفاده از یک مجموعه داده اولیه بدست می‌آورد. این مجموعه داده به دو گروه مثبت و منفی تقسیم‌بندی می‌شوند. قواعد تولید باید گروه مثبت را توصیف کند و هیچ‌یک از موارد موجود در گروه منفی را شامل نشود. بسیاری از عملیاتی که توسط CLIP4 انجام می‌شود، از مسئله SC<sup>8</sup> [4] استفاده می‌کند، به همین دلیل ابتدا توضیحی کلی در خصوص این مسئله خواهیم داد.

### ۲-۱- مسئله SC و کاربرد آن در الگوریتم CLIP4

چند تا از عملیات کلیدی که در الگوریتم CLIP4 انجام می‌شود، از مسئله SC استفاده می‌کند. این مسئله یک نسخه ساده‌تر از مدل برنامه‌نویسی خطی می‌باشد. مدل برنامه‌نویسی خطی، برای بهینه کردن یک تابع با در نظر گرفتن مجموعه‌ای از محدودیت‌ها بکار می‌رود. یکسری ساده‌سازی را اگر در مدل برنامه‌نویسی خطی انجام دهیم به مسئله SC تبدیل خواهد شد. تمام ضرایب تابع هدف و توابع استفاده شده در قسمت محدودیت‌ها مقدار ۱ را دارا می‌باشند، متغیرها از نوع دودویی هستند و تمام توابع محدودیت، بزرگتر مساوی ۱ هستند. مسئله SC یک مسئله NP-Hard است و بنابراین باید از راه‌حل‌های تقریبی برای آن استفاده کرد. شکل ۱ مثالی از یک مسئله SC به صورت یک مدل برنامه‌نویسی خطی را نشان می‌دهد. در الگوریتم CLIP4 از فرم ماتریسی برای نمایش مسئله SC استفاده می‌شود. ماتریسی با عنوان BIN ساخته می‌شود که همانند ماتریس استفاده شده یک ماتریس دودویی بوده و ضرایب توابع محدودیت را مشخص می‌کند. نحوه ساخت آنرا در ادامه بررسی خواهیم کرد. تعداد ستون‌های این ماتریس در حقیقت، مشخص کننده صفات<sup>9</sup> موجود برای رکوردها و به ازای هر نمونه یا رکورد، حاوی یک سطر می‌باشد. در ادامه خواهیم دید که با استفاده از یک نمونه منفی و مقایسه آن با مقادیر موجود در صفات نمونه‌های مثبت منبع داده، اگر مقدار مساوی داشته باشند در سطر و ستون مربوط به آن در این ماتریس ۰ و در غیر اینصورت مقدار ۱ قرار داده می‌شود. به این ترتیب حل این مدل، حداقل ستون‌ها یا صفاتی را برای مجموعه نمونه‌های مثبت موجود، مشخص می‌کند که باعث تمایز بین این نمونه‌ها و نمونه منفی در نظر گرفته می‌شوند. جواب یک بردار است که ستون‌های انتخاب شده حاوی مقدار ۱ می‌باشند.

<sup>7</sup> accurate classification

<sup>8</sup> set covering

<sup>9</sup> attribute

در CLIP4 دلیل استفاده زیاد از این مسئله نیاز به یک راه حل سریع برای آن وجود دارد. در CLIP3 از یک الگوریتم حریصانه<sup>10</sup> ساده برای حل مسئله SC استفاده می شود. در CLIP4 الگوریتمی برای حل این مسئله ارائه شده است که ستون ها را بر اساس کمینه کردن تعداد ۱ هایی که دارا می باشند انتخاب می کند. این الگوریتم باعث افزایش دقت محاسبات و در نتیجه تولید قواعد مناسب تری می شود. شبه کد الگوریتم استفاده شده و مقایسه آن با الگوریتم استفاده شده در CLIP3 در [۱] آمده است.

$$\begin{array}{ll} \text{Minimize :} & \text{Minimize :} \\ x_1 + x_2 + x_3 + x_4 + x_5 = Z & x_1 + x_2 + x_3 + x_4 + x_5 = Z \\ \text{Subject to :} & \text{Subject to :} \\ x_1 + x_3 + x_4 \geq 1 & \begin{bmatrix} 1 & 0 & 1 & 1 & 0 \\ 0 & 1 & 1 & 0 & 1 \\ 0 & 0 & 1 & 1 & 1 \\ 1 & 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \end{bmatrix} \geq 1 \\ x_2 + x_3 + x_5 \geq 1 & \\ x_3 + x_4 + x_5 \geq 1 & \end{array}$$

$$\text{Solution : } Z = 2, \text{ when } x_1 = 1, x_2 = 0, x_3 = 1, x_4 = 0$$

شکل ۱: یک مسئله SC و راه حل آن با استفاده از دو فرم استاندارد و ماتریسی

## ۲-۲- جزئیات الگوریتم CLIP4

شبه کد الگوریتم CLIP4 در شکل ۲ نمایش داده شده است. قبل از شرح مراحل مختلف الگوریتم نمادهای استفاده شده در شبه کد و الگوریتم را در ابتدا تعریف می کنیم. مجموعه نمونه های موجود برای یادگیری، را  $S$  می نامیم. مجموعه نمونه های مثبت را  $S_p$  و مجموعه نمونه های منفی را با  $S_N$  نمایش خواهیم داد. نمونه های مثبت کلاسی را مشخص می کنند که قواعد بر اساس آنها تولید خواهد شد و بقیه را نمونه های منفی می باشند. یک نمونه  $e$  از مجموعه ای از انتخاب گرها به صورت  $s_j = [a_j = v_j]$  تشکیل شده است. الگوریتم CLIP4 قواعدی را به صورت زیر تولید خواهد کرد:

$$IF (s_1 \wedge \dots \wedge s_m) THEN class = class_i \text{ که در آن } s_i = (a_i \neq v_i) \text{ می باشد.}$$

مجموعه های  $S_p$  و  $S_N$  به صورت ماتریس هایی مشخص می شوند. ماتریس نمونه های مثبت با POS نمایش داده شده است و  $N_{POS}$  تعداد نمونه های مثبت موجود یا به عبارت دیگر تعداد سطرهای ماتریس می باشد. در خصوص نمونه های منفی از NEG و برای نمایش تعداد آن از  $N_{NEG}$  استفاده شده است. همانگونه که قبلا اشاره شد، این الگوریتم از ماتریس هایی با عنوان BIN نیز استفاده می کند که دارای  $K$  ستون می باشد، که هر ستونی مشخص کننده یک صفت از نمونه ها می باشد. این ماتریس توسط الگوریتم CLIP4 مقاردهی شده و به صورت یک مسئله SC حل می شود. این الگوریتم در ۳ مرحله کار خود را انجام می دهد که در ادامه آنها را بررسی می کنیم.

### مرحله اول:

داده های موجود در مجموعه نمونه های مثبت به زیرمجموعه هایی با داده های مشابه در یک ساختار درخت تصمیم تبدیل می شود. هر گره از درخت زیرمجموعه ای از داده ها را نشان می دهد. هر سطح درخت با استفاده از یک نمونه منفی ساخته می شود، به این صورت که از آن استفاده می شود تا انتخاب گری که می تواند بین این نمونه منفی و مجموعه نمونه های مثبت تمایز قائل شود را پیدا کنیم. در حین رشد درخت، از هرس کردن استفاده می کنیم تا داده های خطا دار حذف شوند، از رشد زیاد درخت جلوگیری کنیم و زمان اجرا را کاهش دهیم.

<sup>10</sup> greedy algorithm

```

N0=1; Create POS0,1 consisting of entire Sp; Create NEG consisting of entire SN // Initialize
1 for negi, i=1 to NNEG do { // PHASE I STARTS
2   for j=1 to Ni-1 do { //for each POSi-1,j matrix
3     for k=1 to K do { //create new BINj matrix
4       for l=1 to number of POSi-1,j rows do {
5         if posi-1,jl[k] = negi[k] then binjl[k] = 0;
6         if posi-1,jl[k] ≠ negi[k] then binjl[k] = 1;
7         if posi-1,jl[k] = '*' then binjl[k] = 0; // missing value encountered
8         if negi[k] = '*' then binjl[k] = 0; // missing value encountered
9       SOLj = SolveSCProblem(BINj);
10      PruneMatrices(POSi-1,j, SOLj, j=1,...,Ni-1);
11      ApplyGeneticOperators(POSi-1,j, SOLj, j=1,...,Ni-1);
12      Ni=1; //counter for POSi-1 matrices
13     for j=1 to Ni-1 do { //for each POSi,j matrix
14       if POSi-1,j was not pruned or redundant then {
15         for k=1 to K do { //through entire solution vector
16           if SOLj[k]=1 then { //then create new POSi,Ni matrix
17             for l=1 to number of POSi-1,j rows do {
18               if posi-1,jl[k]≠negi[k] then add posi-1,jl to POSi,Ni matrix; }
19             Ni=Ni+1; } }
20       for j=1 to Ni do { //for each POSi matrix check if it large enough to be not considered as a noise
21         if number of rows of POSi,j < NoiseThreshold then { remove POSi,j from the tree; Ni=Ni-1; }
22       EliminateRedundantMatrices(POSi,j j=1,...,Ni); }
23 Create BIN matrix that consist of NNEG columns and NPOS rows, and fill with zeros // PHASE II STARTS
24 for i=1 to NNNEG do { // for all tree leaves
25   for j=1 to NPOS do {
26     for k=1 to number of rows of POSNNEG,i do { if posj0,0 = poskNNEG,i then binj[i]=1; } }
27 SOL = SolveSCProblem (BIN); // select best leaf node subsets
28 for i=1 to NNNEG do { // through entire solution vector
29   create BINi=NEG;
30   if SOL[i]=1 then { // back-project POSNNEG,i matrix
31     for j=1 to K do {
32       for k=1 to NNEG do {
33         if negk[j]='*' then binki[j]=0; else {
34           for l=1 to number of rows of POSNNEG,i do { if poslNNEG,i[j]=negk[j] then binki[j]=0; } } }
35     for j=1 to K do { // convert BINi values to binary
36       for k=1 to NNEG do { if binki[j]≠0 then binki[j]=0; }
37     SOLi = SolveSCProblem (BINi);
38     for j=1 to K do { // start generation of i-th rule
39       if SOLi[j]=1 then { // add selectors to the rule
40         for k=1 to NNEG do { if binki[j]=1 then add "aj≠negk[j]" selector to the Rulei; } } }
// PHASE III STARTS
41 best#covered=0; previous_best#covered=0; // holds # ex. covered by the rule
42 while best#covered≥0 do { //until best rules are accepted
43   for i=1 to NNNEG do { // through all generated rules
44     coversi=NPOS; // # examples covered by Rulei
45     for j=1 to NPOS do { // for all examples in POS
46       for k=1 to K do {
47         for l=1 to number of selectors in Rulei do {
48           if al=k and vk=posj0,1[k] then { coversi=coversi-1; j=j+1; } } } // example not covered by Rulei
49       if best#covered<coversi then best#covered=coversi; best_rule=Rulei; }
50     NPOS = NPOS - best#covered;
51     if NPOS<StopThreshold or NPOS=0 then STOP;
52     POS0,1 = POS0,1 - examples covered by best_rule;
53     if (best#covered<previous_best#covered/2 and best#covered<NPOS/2) then best#covered=-1; // multiple rules
54     previous_best#covered_examples=best#covered_examples; }

```

شكل ٢: شبه كد الگوریتم CLIP4

خطوط ۲-۲۳ شبه کد این مرحله الگوریتم را نشان می‌دهد. درختی که گره‌هایش نمایش دهنده نمونه‌های مثبت می‌باشد با استفاده از یک روش بالا به پایین توسعه داده می‌شود. در سطح  $i$  ام درخت  $N_i$  گره خواهیم داشت که هر کدام با استفاده از ماتریسی با عنوان  $POS_{i,j}(j=1,\dots,N_i)$  نمایش داده می‌شوند. این زیرمجموعه‌ها با استفاده از  $N_{i-1}$  زیرمجموعه سطح بالاتر و یک نمونه منفی  $neg_i$  تولید می‌شوند. هر زیرمجموعه با استفاده از  $neg_i$  به صورت ماتریسی از صفر و یک با عنوان BIN تبدیل می‌شود. که بعداً به صورت یک مسئله SC مدل شده و حل خواهد شد. با استفاده از حل ارائه شده برای مسئله SC، زیرشاخه‌های گره فعلی شکل می‌گیرند و ماتریس‌های POS مربوط به آنها مقداردهی می‌شود. به این ترتیب، از سطح اول تا برگ‌ها با استفاده از نمونه‌های منفی موجود، سطوح مختلف شکل گرفته و مجموعه نمونه‌های مثبت اولیه به زیرمجموعه‌هایی شکسته می‌شوند.

### مرحله دوم:

دو معیار بکار برده می‌شود تا مجموعه‌ای از بهترین زیرمجموعه‌های نهایی (برگ‌های درخت) را انتخاب کنیم. یکی اینکه مجموعه‌های بزرگتر ترجیح داده می‌شوند، زیرا قواعد تولید شده از آنها عمومیت بیشتری خواهد داشت. دوم اینکه به یک معیار کامل بودن<sup>۱۱</sup> نیاز داریم. به این منظور یک عملیات برگشتی برای هر زیرمجموعه موجود در برگ‌ها که حاوی نمونه‌های مثبت می‌باشد با استفاده از تمام نمونه منفی موجود انجام می‌شود، ماتریس بدست آمده به یک ماتریس دودویی تبدیل شده و به صورت یک مسئله SC حل می‌شود. از حل آن برای تولید یک قاعده استفاده می‌شود. اینکار برای تمام زیرمجموعه‌های موجود انجام می‌شود. خطوط ۲۴-۴۱ الگوریتم مربوط به این مرحله می‌باشند. تعداد برگ‌های موجود در درخت  $N_{NEG}$  می‌باشد. از یک ماتریس دودویی BIN استفاده می‌شود تا یک مجموعه کمینه از برگ‌ها انتخاب شود که کلیه نمونه‌های مثبت را پوشش دهد. عملیات برگشتی برای تولید قاعده، بر روی زیرمجموعه‌های انتخاب شده انجام خواهد شد.

### مرحله سوم:

مجموعه قواعد بهینه از میان قواعد موجود انتخاب می‌شود. قواعدی که نمونه‌های مثبت بیشتری را پوشش دهد انتخاب می‌شوند. اگر نتوان از این معیار استفاده کرد، از میان قواعد درگیر در تصمیم‌گیری قاعده‌ای انتخاب می‌شود که طول کمتری داشته باشد، یعنی قاعده‌ای که انتخاب‌گرهای کمتری داشته باشد. خطوط ۴۲-۵۵ الگوریتم شکل ۲ مربوط به این مرحله می‌باشد. برای یافتن تعداد نمونه‌هایی که یک قاعده پوشش می‌دهد به صورت برعکس عمل می‌شود. یعنی تعداد نمونه‌هایی را که پوشش نمی‌دهد و سپس آنرا از تعداد کل نمونه‌های موجود کم می‌کنیم. اینکار به این علت انجام می‌شود که انتخاب‌گرهای موجود از نوع نامساوی هستند. متغیر  $cover_i$  تعداد نمونه‌های مثبت پوشش داده شده توسط قاعده  $i$  ام را نشان می‌دهد. بعد از اینکه یک قاعده مورد قبول قرار گرفت، نمونه‌های متناظر آن از ماتریس POS حذف می‌شوند.

در این الگوریتم، از تعدادی از معیارهای آستانه‌ای<sup>۱۲</sup>، برای هرس کردن، حذف نمونه‌های خطا دار و توقف الگوریتم استفاده می‌شود. این معیارهای آستانه‌ای استفاده شده را در بخش ۵، بررسی خواهیم کرد. نشان داده شده است که پیچیدگی زمانی این الگوریتم  $O(k^3 n^2)$  می‌باشد، که در آن  $k$  تعداد صفات نمونه‌ها و  $n$  تعداد نمونه‌های موجود می‌باشد.

<sup>11</sup> completeness

<sup>12</sup> threshold

### ۳- مثالی از کاربرد الگوریتم CLIP4

در این مثال، ۸ نمونه داریم که در جدول ۱ نشان داده شده‌اند. ۳ نمونه، مقادیر نامشخص دارند که با \* مشخص شده‌اند. کلاس ۱، نمونه‌های مثبت و کلاس ۲، نمونه‌های منفی می‌باشد.

| Ex.# | F1 | F2 | F3 | F4 | Class |
|------|----|----|----|----|-------|
| 1    | 1  | 2  | 3  | *  | 1     |
| 2    | 1  | 3  | 1  | 2  | 1     |
| 3    | *  | 3  | 2  | 5  | 1     |
| 4    | 3  | 3  | 2  | 2  | 1     |
| 5    | 1  | 1  | 1  | 3  | 1     |
| 6    | 3  | 1  | 2  | 5  | 2     |
| 7    | 1  | 2  | 2  | 4  | 2     |
| 8    | 2  | 1  | *  | 3  | 2     |

جدول ۱: نمونه‌های مثبت و منفی بکار برده شده در مثال

بنا به الگوریتم فوق، در ابتدا ماتریس POS برای نمونه‌های مثبت و ماتریس NEG برای نمونه‌های منفی، به صورت زیر تشکیل می‌شود:

$$\begin{array}{c}
 \text{NEG} \\
 \begin{bmatrix} 3 & 1 & 2 & 5 \\ 1 & 2 & 2 & 4 \\ 2 & 1 & * & 3 \end{bmatrix}
 \end{array}
 \quad
 \begin{array}{c}
 \text{POS} \\
 \begin{bmatrix} 1 & 2 & 3 & * \\ 1 & 3 & 1 & 2 \\ * & 3 & 2 & 5 \\ 3 & 3 & 2 & 2 \\ 1 & 1 & 1 & 3 \end{bmatrix}
 \end{array}$$

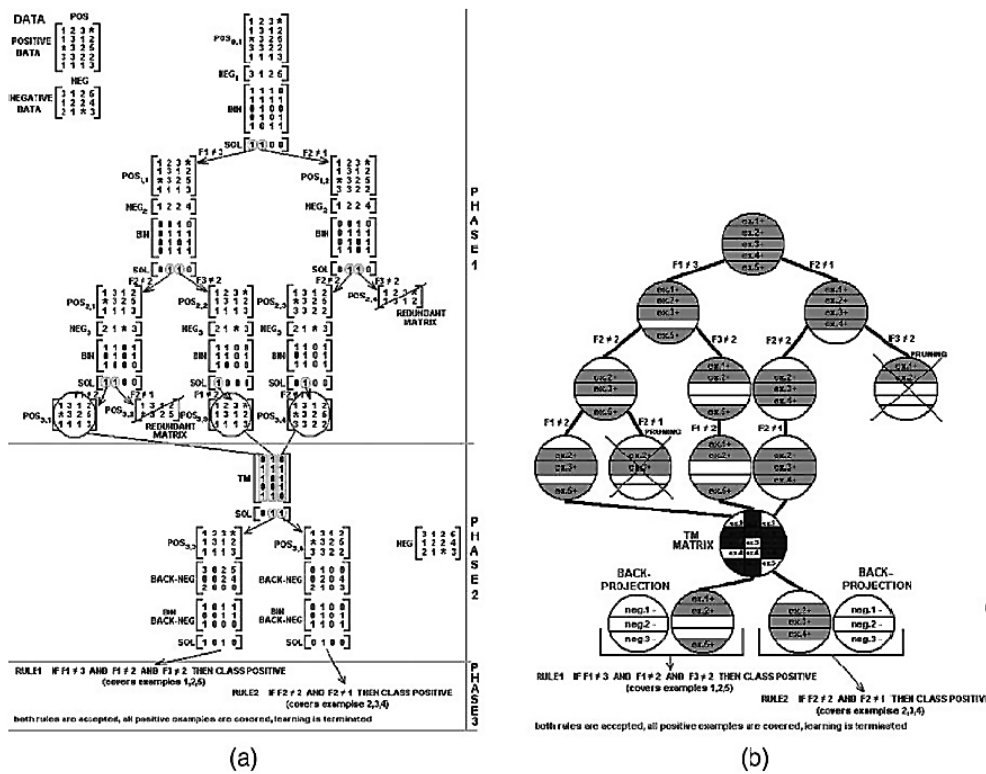
حال برای شروع الگوریتم ماتریس  $POS_{0,1}$  با مقدار ماتریس POS مقارده می‌شود. اولین سطر ماتریس NEG، یعنی نمونه منفی را در نظر می‌گیریم:  $NEG_1 = [3, 1, 2, 5]$ . حال ماتریس BIN به صورتیکه در الگوریتم آمده است با مقایسه مقادیر موجود در نمونه منفی و مقادیر موجود در ماتریس  $POS_{0,1}$  به صورت زیر مقارده می‌شود:

$$\begin{array}{c}
 \text{BIN} \\
 \begin{bmatrix} 1 & 1 & 1 & 0 \\ 1 & 1 & 1 & 1 \\ 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 \end{bmatrix}
 \end{array}$$

این ماتریس به صورت یک مسئله SC، حل می‌شود. جواب عبارتست از:  $SOL = [1, 1, 0, 0]$ . بنابراین با توجه به مقادیر نمونه منفی استفاده شده و SOL بدست آمده، ماتریس  $POS_{0,1}$  با استفاده از دو قاعده  $F1 \neq 3$  و  $F2 \neq 1$  به دو ماتریس  $POS_{1,1}$  و  $POS_{1,2}$  شکسته می‌شود، که به عنوان فرزندان گره فعلی در درخت در نظر گرفته می‌شوند. بقیه مراحل به همین صورت ادامه پیدا می‌کند تا به ازای تمام نمونه‌های منفی سطح درخت ایجاد شود و برگ‌های درخت شکل بگیرد. تولید ماتریس‌های مربوطه در شکل ۳ دیده می‌شود. همانگونه که در شکل دیده می‌شود، برخی از ماتریس‌ها که تعداد نمونه کمی را شامل شده‌اند، حذف می‌شوند. پس از رسیدن به برگ‌ها، مرحله اول الگوریتم به پایان می‌رسد. در مرحله دوم، از میان برگ‌ها، یک مجموعه کمینه انتخاب می‌شود که بتواند کل مجموعه نمونه‌های مثبت را بپوشاند. یک ماتریس BIN دیگر که در اینجا نام آنرا TM گذاشته است، بر اساس رکوردهایی که ماتریس‌های برگ، یعنی  $POS_{3,1}$ ،  $POS_{3,2}$  و  $POS_{3,3}$  از ماتریس POS که حاوی کلیه نمونه‌های مثبت می‌باشد، پوشش می‌دهند، تشکیل می‌شود و سپس

به صورت یک مسئله SC حل می شود و نتیجه آن به صورت  $SOL=[0,1,1]$  بدست می آید که نشان می دهد، دو مجموعه  $POS_{3,3}$  و  $POS_{3,2}$  برای پوشش مجموعه POS کافی می باشند. سپس ماتریس BACK\_NEG برای هر کدام از این دو ماتریس ساخته می شود (مراحل می توانید در شکل ۳ دنبال کنید). برای ساخت این ماتریس ابتدا تمام سلول های آن با ماتریس NEG مقارده می شود، سپس به ازای تمام ستون ها چک می شود که اگر مقداری که در یک سلول قرار دارد، در یکی از رکوردهای ماتریس متناظر POS در همان ستون وجود دارد، مقدار آن به صفر تغییر می کند. همچنین اگر مقدار موجود در ماتریس BACK\_NEG برابر \* یعنی مقدار آن نامشخص باشد، مقدار آن به صفر تغییر می کند. در انتها ماتریس BIN متناظر با جایگذاری مقادیر غیر صفر به ۱، تولید می شود. هر کدام از دو ماتریس BIN تولید شده به صورت یک مسئله SC حل می شود. به ازای هر حل، یک قاعده تولید می شود. قواعد تولید شده به صورت زیر می باشد:

- RULE 1: If  $F1 \neq 3$  and  $F1 \neq 2$  and  $F3 \neq 2$  then class 1  
 RULE 2: If  $F2 \neq 2$  and  $F2 \neq 1$  then class 1



شکل ۳: مثالی از تولید قواعد توسط الگوریتم CLIP4. (a) نمایش سطح پایین اجرا (b) نمایش سطح بالا

#### ۴- ویژگی‌های الگوریتم CLIP4

الگوریتم CLIP4، یکسری از ویژگی‌های الگوریتم CLIP3 را حفظ کرده و ویژگی‌های جدیدی به آنها افزوده است. ویژگی‌های کلی این الگوریتم را در ادامه بررسی می‌کنیم.

##### ۴-۱- رسیدگی به داده‌های نامشخص<sup>۱۳</sup>

داده‌ها یا مقادیر نامشخص، داده‌هایی هستند که در رکوردهای موجود در بانک‌اطلاعاتی وارد نشده‌اند. الگوریتم CLIP4، رکوردهایی را که مقادیر نامشخص دارند، را حذف نمی‌کند و همچنین آنها را با مقادیر پیش‌فرض یا محاسبه‌شده مقداردهی نمی‌کند. این الگوریتم از مقادیر موجود برای بقیه صفات رکورد استفاده کرده و از مقادیر نامشخص صرف‌نظر می‌کند. بنابراین در این الگوریتم کاربر نیازی به انجام پردازش بر روی داده‌های نامشخص موجود در منبع‌داده قبل از اعمال الگوریتم نخواهد داشت.

##### ۴-۲- رسیدگی به صفات اسمی<sup>۱۴</sup> و تولید قاعده برای منابع داده با چند کلاس

الگوریتم CLIP4 داده‌های اسمی را در یک مرحله پیش‌پردازش به صورت عددی کدگذاری می‌کند. الگوریتم وابسته به شمای کدگذاری استفاده شده نیست، زیرا در هنگام تولید قواعد از معیاری که مرتبط با محاسبه فاصله باشد، استفاده نمی‌کند. همانگونه که دیدیم، این الگوریتم از مسئله SC استفاده می‌کند، که وابسته به مقادیر واقعی صفات نمی‌باشد. برای مواردی که چندکلاس داریم، الگوریتم در مراحل مختلف و در هر مرحله با در نظر گرفتن یکی از کلاس‌ها به عنوان کلاس مثبت، قواعد مربوط به آنرا تولید می‌کند.

##### ۴-۳- معیارهای آستانه‌ای

در طول الگوریتم از استفاده شده است که برای هرس کردن و حذف نمونه‌های خطا دار از آن استفاده می‌کند. آستانه‌های استفاده شده به صورت زیر می‌باشد:

آستانه خطا<sup>۱۵</sup> (NT) مشخص می‌کند که چه گره‌هایی باید حذف شوند. این آستانه حداقل داده‌های موجود در یک زیرمجموعه را تعیین می‌کند. گره‌ای که زیرمجموعه مربوطه آن دارای تعداد نمونه‌های کمتر از این آستانه باشد، حذف خواهد شد.

آستانه هرس<sup>۱۶</sup> (PT) نیز همانند آستانه خطا برای حذف گره‌ها از درخت مورد استفاده قرار می‌گیرد. الگوریتم از یک تابع کمی برای اندازه‌گیری میزان شایستگی<sup>۱۷</sup> یک گره استفاده می‌کند. سپس تعداد PT گره به ترتیب میزان شایستگی آنها انتخاب می‌شوند.

آستانه توقف<sup>۱۸</sup> (ST) زمان توقف الگوریتم را مشخص می‌کند. الگوریتم وقتی متوقف می‌شود که تعداد نمونه مثبت کمتر از ST پوشش نداده شده توسط قواعد باقی بماند. همانگونه که قبلاً گفته شد، وقتی قاعده‌ای انتخاب شد، نمونه‌های پوشش داده شده توسط آن قاعده از ماتریس POS که حاوی نمونه‌های مثبت موجود می‌باشد، حذف می‌شوند.

<sup>13</sup> missing values

<sup>14</sup> nominal

<sup>15</sup> noise threshold

<sup>16</sup> pruning threshold

<sup>17</sup> fitness

<sup>18</sup> stop threshold

## ۵- نتایج تجربی

الگوریتم CLIP4 آنگونه که در [1] آمده است، بر روی چند مسئله محک‌زنی تست شده است. تست‌های انجام شده برای مقایسه این الگوریتم با الگوریتم‌های موجود و همچنین تست‌هایی برای مشخص کردن ویژگی‌های خاص الگوریتم می‌باشد.

### ۵-۱- منابع داده

CLIP4 بر روی ۲۹ منبع داده تست شده است. مشخصات این منابع داده را به صورت زیر می‌توان دسته بندی کرد:

- اندازه منابع داده یادگیری<sup>۱۹</sup>: بین ۱۵۰ تا ۳۰ هزار نمونه
- اندازه منابع داده تست: بین ۴۵ تا ۵۶۵ هزار نمونه
- تعداد صفات: بین ۴ تا ۱۵۵۸
- تعداد کلاس‌ها: بین ۲ تا ۱۰

منابع داده از طریق مخزن منابع اطلاعاتی موجود در [6] و [7] بدست آمده است. اطلاعات کلی این منابع داده استفاده شده که در جدول ۲، که با حروف اختصار نمایش داده شده است، در [1] آورده شده است.

### ۵-۲- مقایسه با سایر الگوریتم‌ها

نرخ خطا، تعداد قواعد، تعداد انتخاب‌گرها و زمان اجرای CPU برای الگوریتم CLIP4 با استفاده از ۱۶ منبع داده در جدول ۲، دیده می‌شود. نتایج بدست آمده از الگوریتم CLIP4 با نتایج بدست آمده در [8] مقایسه شده است. در [8] کارایی ۳۳ الگوریتم یادگیری ماشین با یکدیگر مقایسه شده‌اند. در جدول ۲، تنها مقادیر بیشینه و کمینه برای موارد تستی مختلف در کل ۳۳ الگوریتم تست شده آورده شده است. همانگونه که دیده می‌شود، میانگین نرخ خطا برای الگوریتم CLIP4 عدد ۲۵,۱ درصد بدست آمده است. کمترین میزان میانگین نرخ خطا مربوط به الگوریتم POLYCLASS [9] می‌باشد که مقدار ۱۹,۵ درصد را بدست آورده است. در [8] معیاری با عنوان اهمیت آماری برای نرخ‌های خطا در نظر گرفته شده است و نشان داده شده است که تفاوت میانگین نرخ خطای دو الگوریتم در یک سطح ۱۰ درصد در صورتیکه از ۵,۹ تجاوز کند، اهمیت دارد. بنابراین الگوریتم CLIP4 نیز در گروه مشابه همان الگوریتم قرار می‌گیرد و به این صورت یکی از الگوریتم‌های یادگیری ماشین خوب به حساب می‌آید.

زمان اجرای CPU در الگوریتم CLIP4 برای ۱۶ منبع داده، در جدول ۲، دیده می‌شود. میانگین زمان اجرای الگوریتم، ۵,۸ دقیقه می‌باشد. نتایج [8] نشان می‌دهد که میانگین زمان اجرا برای ۳۳ الگوریتم تست شده بین ۶,۹ ثانیه و ۴۶,۸ ساعت بوده است. که آنها را به دو دسته کمتر از ۱۰ دقیقه (۲۲ الگوریتم) و بیشتر از ۱۰ دقیقه (۱۱ الگوریتم) تقسیم‌بندی کرده است. این الگوریتم در دسته اول قرار می‌گیرد.

میانگین تعداد قواعد تولید شده توسط CLIP4، ۱۶,۸ می‌باشد. [8] میانگین تعداد برگ‌های درخت تولید شده توسط ۲۱ الگوریتم درخت تصمیم تست شده (که برابر تعداد قواعد تولید شده می‌باشد)، برابر ۱۷,۸ می‌باشد. همانگونه که دیده می‌شود، الگوریتم CLIP4، از لحاظ میانگین قواعد تولید شده توسط آن میزان کمتری را دارا می‌باشد. به طور کلی، ۱۰ الگوریتم از ۲۱ الگوریتم‌های تست شده در [8]، میانگین قواعد تولید شده توسط آنها کمتر از الگوریتم CLIP4 می‌باشد. میانگین تعداد انتخاب‌گرهای تولید شده توسط الگوریتم CLIP4 برابر ۱۶ منبع داده برابر ۵۸۹,۵ می‌باشد که به ازای هر قاعده تعداد ۳۵ انتخاب‌گر بدست می‌آید.

<sup>19</sup> training dataset

| تعداد انتخاب‌گرهای<br>CQIP4 قواعد | زمان CPU<br>الگوریتم<br>CQIP4 (ثانیه) | زمان CPU گزارش شده برای ۳۳<br>الگوریتم یادگیری ماشین |           | تعداد قواعد<br>تولید شده<br>توسط<br>CQIP4 | میانگین<br>گزارش شده<br>برای تعداد<br>قواعد تولید<br>شده توسط<br>۲۱<br>الگوریتم | نرخ<br>خطای<br>CQIP4 | نرخ‌های خطای گزارش<br>شده برای الگوریتم<br>یادگیری ماشین |       | منبع داده |
|-----------------------------------|---------------------------------------|--|-----------|---|---|----------------------|--|-------|-----------|
|                                   |                                       | بیشینه (h)   | کمینه (s) |   |   |                      | بیشینه   | کمینه |           |
| 121.6                             | 5.1                                   | 2.7  | 4         | 4.2                                       | 7   | 5                    | 9  | 3     | bcw       |
| 60.7                              | 46                                    | 23.9   | 12        | 8   | 15  | 53                   | 60   | 43    | cmc       |
| 90                                | 662.8                                 | 475.2  | 2         | 8   | 13  | 9                    | 38   | 5     | dna       |
| 192.3                             | 2.2                                   | 3.3  | 4         | 11.6                                      | 6   | 28                   | 34   | 14    | hea       |
| 133.5                             | 35.8                                  | 5.5  | 9         | 10.5                                      | 11  | 29                   | 31   | 22    | bos       |
| 189                               | 166.4                                 | 12.4   | 1         | 41  | 24  | 29                   | 82   | 27    | led       |
| 272.4                             | 6.6                                   | 1.5  | 5         | 9.7                                       | 10  | 37                   | 43   | 28    | bld       |
| 64.1                              | 5.8                                   | 2.5  | 7         | 4   | 7   | 29                   | 31   | 22    | pid       |
| 3199                              | 3696.7                                | 73.2   | 8         | 61  | 63  | 20                   | 40   | 10    | sat       |
| 1169.9                            | 614.6                                 | 75.6   | 28        | 39.2                                      | 39  | 14                   | 52   | 2     | seg       |
| 242                               | 90.5                                  | 3.8  | 1         | 18  | 2   | 32                   | 45   | 30    | sno       |
| 119                               | 164.2                                 | 16.1   | 3         | 4   | 12  | 1                    | 89   | 1     | thy       |
| 380.7                             | 45.3                                  | 14.1   | 14        | 21.3                                      | 38  | 44                   | 49   | 15    | veh       |
| 51.7                              | 4.4                                   | 25.2   | 2         | 9.7                                       | 2   | 6                    | 6  | 4     | vot       |
| 85                                | 43.1                                  | 4.3  | 4         | 9   | 16  | 25                   | 48   | 15    | wav       |
| 273.2                             | 0.7                                   | 10.2   | 6         | 9.3                                       | 20  | 40                   | 69   | 33    | tae       |
| ۵۹۸,۵                             | ۵,۸ دقیقه                             | ۴۶,۸   | ۶,۹       | ۱۶,۸                                      | ۱۷,۸  | ۲۵,۱                 | ۴۵,۴   | ۱۷,۱  | میانگین   |

جدول ۲: نتایج مقایسه با سایر الگوریتم‌ها

### ۵-۳- تست قابلیت انعطاف<sup>۲۰</sup>

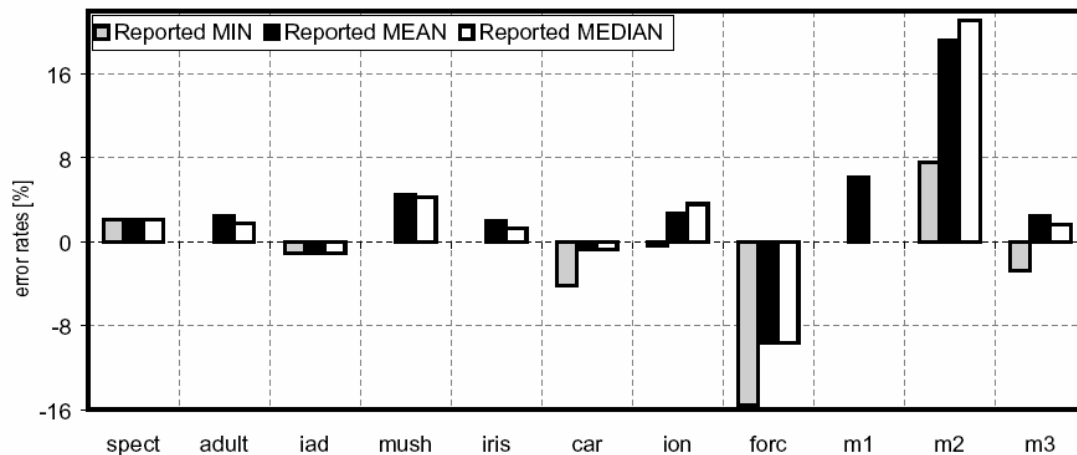
تست‌های انجام شده، شامل اندازه‌گیری نرخ‌خطا، تعداد قواعد تولید شده، تعداد انتخاب‌گرها و زمان اجرای CPU برای ۱۱ منبع داده می‌باشد. تست‌ها نشان می‌دهند که بسته به نوع معیارهای آستانه‌ای انتخاب شده، نتایج متفاوتی خواهیم داشت. جدول مربوط به نتایج در [1] آمده است. نتایج نشان می‌دهد که در منبع داده adult، نرخ خطا ۳,۱ درصد افزایش خواهد یافت، اگر بخواهیم که تعداد قواعد تولید شده را ۳,۵ برابر کاهش دهیم. و همین افزایش نرخ خطا در ازای تولید تعداد قواعد کمتر برای دیگر منابع داده نیز تکرار شده است.

قابلیت انعطاف الگوریتم با بکارگیری آن برای منبع داده adult که تعداد نمونه (رکورد) زیادی دارد و همچنین منبع داده iad که تعداد صفات زیادی دارد، تست شده است. نرخ خطا برای منبع داده adult مشابه کمترین نرخ گزارش شده می‌باشد، که مربوط به الگوریتم FSS Naïve Bayes می‌باشد. نرخ خطا برای منبع داده iad برابر ۳,۹ درصد می‌باشد. بهترین نرخ خطای گزارش شده برای این منبع داده، الگوریتم C4.5 می‌باشد که نرخ خطای آن ۲,۸ درصد می‌باشد. ولی این در حالی است که تعداد قواعد تولید شده توسط CLIP4 برابر ۶,۷ قاعده می‌باشد، که ۳,۷ برابر کمتر از تعداد قواعد تولید شده توسط C4.5 می‌باشد. همچنین الگوریتم CLIP4 برای این منبع داده از ۳۰ انتخاب‌گر در هر قاعده استفاده می‌کند. لازم به ذکر است که این منبع داده دارای ۱۵۵۸ صفت می‌باشد. بنابراین می‌توان گفت که الگوریتم CLIP4 در

<sup>20</sup> flexibility

مورد منابع داده با تعداد صفات زیاد، به خوبی عمل می‌کند و با تولید تعداد قواعد کمتر، تقریباً نرخ خطای زیادی را نیز تولید در بر ندارد.

شکل ۴، تفاوت نرخ خطای بدست آمده توسط الگوریتم CLIP4، در مقایسه با نتایج گزارش شده در مقالات برای ۱۱ منبع داده نشان داده شده است. مقادیر مثبت، یعنی CLIP4 بهتر از بقیه الگوریتم‌های گزارش شده عمل کرده است.



شکل ۴: مقایسه بین نرخ خطاهای گزارش شده و نرخ خطای CLIP4 (مقادیر مثبت بهبود را نشان می‌دهد)

## ۶- جمع‌بندی و نتیجه‌گیری

در این مقاله جزئیات الگوریتم CLIP4 را بررسی کردیم. این الگوریتم یکی از الگوریتم‌های یادگیری ماشین استقرایی به حساب می‌آید که برای دسته‌بندی<sup>۲۱</sup> مورد استفاده قرار می‌گیرد. این الگوریتم از یک ایده ترکیبی از الگوریتم‌های مبتنی بر قاعده و الگوریتم‌های درخت تصمیم، استفاده می‌کند. از ویژگی‌های بارز آن می‌توان تولید قواعدی که حاوی نامساوی هستند و رسیدگی به منابع داده با ابعاد زیاد را ذکر کرد. تولید قواعد نامساوی این امکان را می‌دهد که برای یک مفهوم خاص به تعداد قواعد کمتری نیاز داشته باشیم، و همانگونه که قسمت بررسی نتایج آزمایشات انجام شده بر روی این الگوریتم مشاهده شد، در برخی از موارد می‌تواند با تولید تعداد قواعد کمتر، نرخ خطای قابل قبولی را نیز به همراه داشته باشد. این الگوریتم علاوه بر دسته‌بندی در موارد دیگری مانند تشخیص مهمترین صفات موجود در مجموعه داده نیز کاربرد دارد. همانگونه که در قسمت مقدمه ذکر شد، چهار فاکتور مهمی که یک الگوریتم یادگیری ماشین که برای دسته‌بندی بکار می‌رود، باید دارا باشد عبارتند از: دقت دسته‌بندی دقیق، قواعد ساده، تولید کارای قواعد و انعطاف‌پذیری می‌باشد. نشان داده شد که اجرای این الگوریتم بر روی منابع داده مختلف، در حضور داده‌های خطادار و نامشخص امکان‌پذیر است و نرخ خطای بدست آمده از لحاظ آماری شبیه به بهترین الگوریتم گزارش شده می‌باشد. همچنین استفاده از هرس کردن‌های پی‌درپی در طول اجرا و مسئله SC، به الگوریتم توانایی تولید یک مجموعه کمینه از قواعد مختصر را می‌دهد. همچنین زمان اجرای مورد نیاز الگوریتم بدلیل استفاده از هرس کردن به میزان قابل قبولی کاهش یافته است و همانگونه که دیده شد، در بین ۲۲ الگوریتم از ۳۳ الگوریتم موجود قرار گرفت. همچنین با آزمایش‌های انجام شده، مشخص شد که الگوریتم از قابلیت انعطاف بالایی برخوردار است و می‌تواند رسیدگی به منابع داده با نمونه‌های زیاد و همچنین مجموعه صفات زیاد را، به صورت قابل قبولی انجام دهد.

<sup>21</sup> classification

## مراجع:

- [1] K.J. Cios, L.A. Kurgan, CLIP4: Hybrid inductive machine learning algorithm that generates inequality rules, *Information Sciences* 163 (2004) 37–83.
- [2] K.J. Cios, N. Liu, An Algorithm which Learns Multiple Covers via Integer Linear Programming. Part I—The CLILP2 Algorithm, *Kybernetes*, 24:2, pp. 29–50 (The Norbert Wiener Outstanding Paper Award), 1995.
- [3] K.J. Cios, D.K. Wedding, N. Liu, CLIP3: Cover learning using integer programming, *Kybernetes* 26 (4–5) (1997) 513–536.
- [4] M. Bellmore, H.D. Ratliff, Set covering and involuntary bases, *Management Science* 18 (1971)3.
- [5] P. Clark, R. Boswell, Rule induction with CN2: Some recent improvements, in: Y. Kodratoff (Ed.), *Lecture Notes in Artificial Intelligence*, Proceedings of European Working Session on Learning, Springer-Verlag, 1991, pp. 151–163.
- [6] C.L. Blake, C.J. Merz, UCI Repository of Machine Learning Databases. Available from <<http://www.ics.uci.edu/~mllearn/MLRepository.html>>. Irvine, CA: University of California, Department of Information and Computer Science, 1998.
- [7] P. Vlachos, StatLib Project Repository. Available from <<http://lib.stat.cmu.edu>>, 2000.
- [8] T.-S. Lim, W-Y. Loh, Y-S. Shih, A comparison of prediction accuracy, complexity, and training time of thirty-three old and new classification algorithms, *Machine Learning* 40 (2000) 203–228.
- [9] C. Kooperberg, S. Bose, C.J. Stone, Polychotomous regression, *Journal of the American Statistical Association* 92 (1997) 117–127.