

آداده کاوی

dataacademy.ir

روش جدیدی برای داده کاوی مبتنی بر مجموعه های دانه درشت و پردازش تکاملی

^۱بابک نصیری

^۱مهران محسن زاده

^۱عبدالرضا سالاری

- ۱-دانشگاه آزاد اسلامی ، واحد علوم و تحقیقات

دانشکده فنی - گروه کامپیوتر

abishSalar@yahoo.com

m_mohsenzadeh77@yahoo.com

nasiri_babak@yahoo.com

^۲محمد تشه لب

- ۲-دانشگاه خواجه نصیرالدین طوسی، دانشکده برق

teshnehlab@eetd.kntu.ac.ir

چکیده :

درخت تصمیم یکی از روش‌های طبقه‌بندی اطلاعات در داده کاوی می‌باشد که به دلیل سادگی و قابلیت تفهیم ، بسیار مورد استفاده قرار می‌گیرد . در این مقاله روش جدیدی برای تولید درخت تصمیم^۱ مبتنی بر پردازش تکاملی و مجموعه دانه درشت ارائه شده است . ابتدا بررسی مجموعه داده‌ها ، عملیات پیش پردازش شامل استفاده از مقادیر عمومی به جای مقادیر گمشده^۲ و گستره سازی مشخصه‌ها^۳ انجام ، سپس با بهره گیری از تئوری مجموعه های دانه درشت^۴ و پردازش تکاملی ، مشخصه‌های مناسب استخراج می‌شوند و در انتها با استفاده از پردازش تکاملی درخت تصمیم ساخته می‌شود . نتایج بدست آمده بر روی مجموعه داده های مختلف که دارای مشخصه های زیادی هستند ، در مقایسه با نتایج معترض ارائه شده ، بسیار مطلوب بوده است.

کلمات کلیدی : داده کاوی ، درخت تصمیم ، مجموعه دانه درشت ، پردازش تکاملی ، استخراج مشخصه ها^۵

۱- مقدمه

داده کاوی یکی از مهمترین مراحل در فرآیند اکتشاف دانش بشمار می‌رود و وظیفه آن استخراج دانش و الگوهایی است که بصورت ضمنی در داده‌ها وجود دارد . داده کاوی بر روی مجموعه داده‌های بزرگ هزینه محاسباتی زیادی دارد ، به همین علت نیاز است که مشخصه‌های اصلی داده‌ها از بین کل مشخصه‌ها ، قبل از عمل داده کاوی استخراج شود .

¹ Decision Tree

² Missing Values

³ Feature Discretization

⁴ Rough Set Theory

⁵ Feature Selection

در دهه گذشته روش‌های مختلفی برای استخراج مشخصه‌ها پیشنهاد شده است که به طور کلی می‌توان آنها را به دو گروه زیر تقسیم نمود. الگوریتم‌های رتبه‌بندی مشخصه‌ها و الگوریتم‌های یافتن زیرمجموعه مینیمم. در روش رتبه‌بندی مشخصه‌ها، به هر مشخصه یک رتبه نسبت می‌دهند. در این روش زیرمجموعه مینیمم از مشخصه‌ها برای آنالیز‌های بعدی مشخص نمی‌شود و فقط فاصله بین مشخصه‌ها بدست می‌آید اما در روش یافتن زیرمجموعه مینیمم، این موضوع بصورت معکوس وجود دارد. [10]

برای تولید درخت تصمیم نیز روش‌های مختلفی پیشنهاد شده است که از بین آنها می‌توان به روش‌های ID3، C4.5 و پردازش تکاملی اشاره نمود.

در این مقاله روش جدیدی برای داده کاوی با استفاده از درخت تصمیم مبتنی بر پردازش تکاملی ارائه شده است. همچنین از پردازش تکاملی برای استخراج مشخصه‌های داده‌ها از داده‌های حجمی استفاده گردیده است. روش پیشنهادی پیاده‌سازی و نتایج اجرا با سایر روش‌ها مقایسه شده است.

مجموعه داده‌های استفاده شده برای ارزیابی روش برگرفته از کار گروه‌های تحقیقاتی بر روی داده کاوی و وب کاوی است که در گردهمایی COIL 2000 [3] به مسابقه گذاشته شده است. نتایج بدست آمده از اجرای روش ارائه شده در این مقاله نسبت به نتایج روش‌های مشابه بسیار مطلوب بوده است.

روش ارائه شده دارای سه قدم اصلی است. پیش پردازش، استخراج مشخصه‌ها و ایجاد درخت تصمیم. هر یک از این قدم‌ها تشریح شده است. در بخش دوم مقاله عملیات مورد نیاز در پیش پردازش داده‌ها مورد بررسی قرار گرفته است. در بخش سوم روشی برای بهبود استخراج مشخصه‌ها، با استفاده از تئوری مجموعه‌های دانه درشت و پردازش تکاملی مورد استفاده قرار گرفته است. در بخش چهارم روشی برای تولید درخت تصمیم با استفاده از پردازش تکاملی توضیح داده شده است. در بخش پنجم نتایج حاصل از پیاده‌سازی روش پیشنهادی بر روی مجموعه داده‌های مختلف مورد ارزیابی قرار گرفته است و بخش آخر به نتیجه گیری و پیشنهادهایی برای بهبود تخصیص دارد.

۲- پیش پردازش

در این بخش پیش پردازش‌های لازم بر روی داده‌ها، شامل عملیاتی برای پاکسازی داده‌ها^۶ نظر استفاده از مقدار عمومی به جای مقادیر گمشده و گسته سازی مشخصه‌ها تشریح شده است.

روش‌های مختلفی برای مدیریت نمودن مقادیر گمشده وجود دارد که می‌توان به حذف رکورد مربوطه (در صورتیکه مشخصه گمشده، مشخصه کلاس باشد)، پر کردن مقدار برای مشخصه گمشده به صورت دستی، جایگزین نمودن یک مقدار عمومی به جای مقدار مشخصه گمشده، جایگزین نمودن میانگین مقادیر موجود در ستون مشخصه گمشده به جای مقدار مشخصه گمشده اشاره نمود. در روش پیشنهادی از جایگزین نمودن یک مقدار عمومی به جای مقدار مشخصه گمشده استفاده شده است. [10]

موضوع بعدی گسته سازی مشخصه‌ها - چگونگی تبدیل داده‌های با بازه بزرگ به داده‌های با بازه کوچکter است. برای حل این مسئله در روش پیشنهادی از الگوریتم ChiMerge [2] که از توزیع آماری χ^2 برای گسته سازی بازه‌های عددی استفاده می‌کند، استفاده شده است. این الگوریتم شامل یک مرحله آغازین و یک فرآیند پائین به بالای اتصال می‌باشد، که بازه‌ها به طور پیوسته تا هنگامی به هم متصل می‌شوند که شرط پایانی برقرار گردد. الگوریتم فوق در ابتدا با مرتب

⁶ Data Cleansing

سازی داده های آموزشی که شامل مقادیری است که باید گسته شوند، شروع به کار می کند و اولین مجموعه گسته شده با فرار گرفتن هر مقدار در بازه خود شکل می گیرد. فرمول محاسبه مقدار χ^2 به شرح زیر می باشد:

$$\chi^2 = \sum_{i=1}^m \sum_{j=1}^k \frac{(A_{ij} - E_{ij})^2}{E_{ij}}$$

که در آن

$m = 2$ (دو بازه با هم مقایسه می شوند)

$k =$ تعداد کلاس ها

$A_{ij} =$ تعداد نمونه ها در i امین بازه j امین کلاس

$\sum_{j=1}^k A_{ij} = R_i$ = تعداد نمونه ها در i امین بازه

$\sum_{j=1}^m A_{ij} = C_i$ = تعداد نمونه ها در j امین کلاس

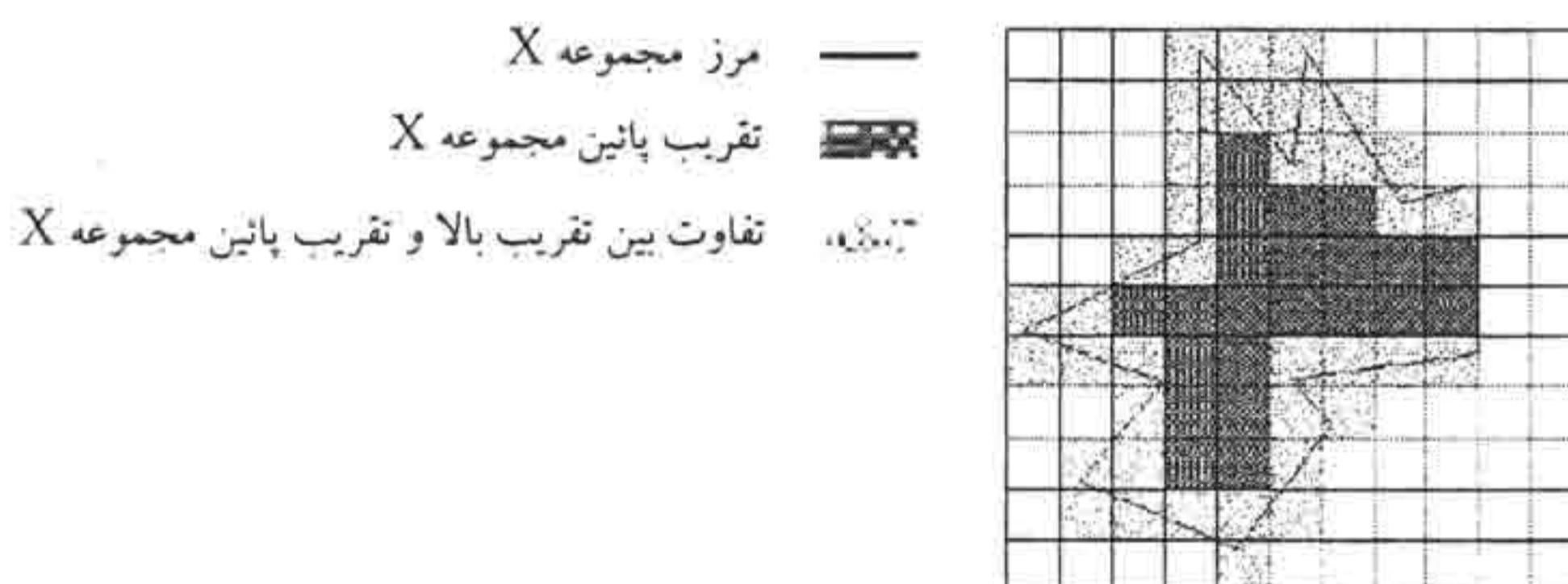
$A_{ij} \sum_{j=1}^k C_j = N$ = مجموع تعداد نمونه ها

$A_{ij} = \frac{R_i \times C_j}{N}$ = فرکانس مورد نظر از

می بلشند.

۳- استخراج مشخصه ها با استفاده از مجموعه های دانه درشت و پردازش تکاملی

تئوری مجموعه های دانه درشت بر پایه تقریب مفاهیم (مجموعه) با استفاده از دانش موجود در یک سیستم اطلاعاتی می باشد. از نقاط قوت این تئوری می توان به کارایی آن در مورد مقادیر گمشده اشاره نمود. تئوری مجموعه های دانه درشت برای پردازش داده های غیر کامل، دو مجموعه تقریب بالا و تقریب پائین را تعریف می کند.



شکل (۱): شماتیکی مجموعه های دانه درشت.

چنانچه X را مجموعه واقعی در نظر بگیریم، تقریب بالا شامل اشیائی است که احتمالاً عضو X می باشد، در حالیکه تقریب پائین شامل اشیائی است که قطعاً عضو X هستند. در صورتیکه تقریب بالا و پائین برای یک زیرمجموعه از سیستم اطلاعاتی برابر باشند آنگاه آن زیرمجموعه را قابل تعریف و در غیر این صورت آن مجموعه را بسختی^۷ قابل تعریف می نامند.

⁷ Rough

سیستم اطلاعاتی عبارت از یک زوج مرتب $(S = (U, A \cup \{d\})$ که U یک مجموعه متناهی و غیر تهی (به نام مجموعه جهانی) و A یک مجموعه متناهی و غیر تهی از مشخصه های شرطی و d مشخصه تصمیم می باشد [4,5]

ماتریس تمایز نیز یک ماتریس مستقران با ابعاد $|U| \times |U|$ با درایه های C_{ij} می باشد که بصورت زیر تعریف می شود:

$$C_{ij} = \{a \in A \mid a(x_i) \neq a(x_j)\} \text{ for } i, j = 1, \dots, n\}$$

تابع تمایز نیز با ترکیب فصلی کلیه مشخصه های موجود در \mathbb{Z}_j و سپس ترکیب عطفی کلیه \mathbb{Z}_i ها با یکدیگر ساخته می شود:

$$f_d(a_1, \dots, a_m) = \wedge \{\vee c_{ij} \mid 1 \leq i \leq n, c_{ij} \neq 0\}$$

برای استخراج مشخصه های مناسب، با استفاده از تئوری مجموعه های دانه درشت، لازم است که ابتدا ماتریس تمایز ساخته شود و سپس تابع تمایز از روی این ماتریس بدست آید. با ساده سازی این تابع، مجموعه کاهش یافته (شامل مشخصه های مناسب) بدست می آید. هنگامیکه حجم داده ها زیاد باشد، ساده سازی تابع تمایز بدست آمده دارای هزینه محاسباتی بسیار بالایی می باشد. از اینرو بجای ساده سازی تابع تمایز، از توانایی الگوریتم های ژنتیکی که جستجوی بهینه در فضاهای بزرگ است، استفاده می شود.

ساختار کروموزوم ها بصورت یک رشته دودویی در نظر گرفته شده و تعداد ژن ها برابر با تعداد مشخصه ها در مجموعه داده می باشد. برای مثال، اگر اشیاء یک سیستم اطلاعاتی دارای 6 مشخصه (با مشخصه تصمیم) باشد، کروموزوم مورد استفاده توسط این روش شامل 6 ژن می باشد.

تابع تطبیق استفاده شده در این روش، برای انتخاب بهترین ها، از ماتریس تمایز استفاده می کند. با فرض اینکه $b(i, (k, n))$ یک درایه از ماتریس تمایز مربوط به مشخصه i ام و دو تایی (O_k, O_n) باشد، ماتریس تمایز بصورت زیر ساخته می شود [6,7]:

$$\text{for } i \in \{1, \dots, N\} : b(i, (k, n)) = \begin{cases} 1 & a_i(O_k) \neq a_i(O_n) \\ 0 & a_i(O_k) = a_i(O_n) \end{cases}$$

$$b(N+1, (k, n)) = \begin{cases} 0 & d_i(O_k) \neq d_i(O_n) \\ 1 & d_i(O_k) = d_i(O_n) \end{cases}$$

تابع تطبیق نیز بصورت زیر می باشد که در آن بخش اول تعداد یک های موجود در کروموزوم $(N_r)_r$ و بخش دوم شامل تعداد سطرهای پوشش داده شده از ماتریس تمایز B_r می باشد.

$$F(r) = (1 - \alpha) * \text{Min}(\varepsilon_1, \frac{N - N_r}{N}) + \alpha * \text{Min}(\varepsilon_2, \frac{C_r}{K}) \quad \text{where}$$

$$K = \frac{m^2 - m}{2}$$

همچنین N تعداد مشخصه های اولیه، m تعداد اشیاء در سیستم اطلاعاتی، K تعداد سطرهای ماتریس تمایز، ε_1 حداقل درصد مشخصه های حذف شده قابل قبول، ε_2 حداقل درصد پوشش های داده شده مورد قبول و α میزان اهمیت به هر یک از دو بخش تابع می باشد.

۴- ساخت درخت تصمیم با استفاده از پردازش تکاملی

الگوریتم های متعددی برای ساخت درخت تصمیم همانند $C4.5$ ، $ID3$ وجود دارند. اینگونه الگوریتم ها هنگامیکه تعداد مشخصه ها زیاد باشند هزینه محاسباتی بالایی دارند. یکی از راه های ساخت درخت تصمیم استفاده از الگوریتم ژنتیک می باشد.

در روش پیشنهادی برای ساخت درخت تصمیم از الگوریتم های ژنتیک استفاده می شود. برای این منظور از الگوریتم ارائه شده در [9] استفاده شده است. برای محاسبه تابع تطابق، هر درخت تولید شده توسط الگوریتم، به طبقه بندی مجموعه داده های آزمایشی می پردازد. بدین ترتیب هر درختی که بتواند تعداد بیشتری از داده های آموزشی را بطور درست طبقه بندی کند، دارای مقدار تابع تطابق بالاتری خواهد بود.

در ابتدا تعدادی درخت تصادفی تولید می شود، سپس درخت هایی که دارای مقدار تابع تطابق بالاتری هستند، انتخاب شده و با یکدیگر ترکیب می شوند، تا نتیجه بهتری حاصل شود. همچنین برای همبری از یک نوع خاص استفاده شده است بدین ترتیب که ابتدا یک نقطه قطع انتخاب می شود و سپس زیر درخت از آن گره مشترک در دو درخت با یکدیگر تعویض می شوند (در صورتیکه در زیر درخت ها گره تکراری موجود نباشد).

۵- پیاده سازی و ارزیابی روش پیشنهادی

در این بخش از مقاله نتایج حاصل از پیاده سازی روش پیشنهادی بر روی مجموعه داده های مختلف ارائه و مورد ارزیابی و مقایسه با سایر روشها قرار گرفته است.

مجموعه داده های انتخاب شده برای ارزیابی روش پیشنهادی شامل دو مجموعه *Adeater* و *Insurance* می باشد. مجموعه داده *Adeater* برگرفته از مشخصه های بکار گرفته شده در حذف عکس های تبلیغاتی در صفحات اینترنت می باشد [5]. این مجموعه داده دارای ۳ مشخصه عددی و ۱۵۵۵ مشخصه دودویی است. مجموعه داده *Insurance* نیز دارای ۸۵ مشخصه عددی می باشد، که این داده ها برگرفته شده از اطلاعات یمه گذاران می باشد و برای پیش بینی سرمایه گذاری در یک نوع بیمه خاص مورد استفاده قرار گرفته است. هر دو مجموعه مذکور همچنین، در گردهمایی COIL 2000 [3] به مسابقه گذاشته شده است. در آن مسابقات روش های مختلف داده کاوی مورد ارزیابی قرار می گیرند.

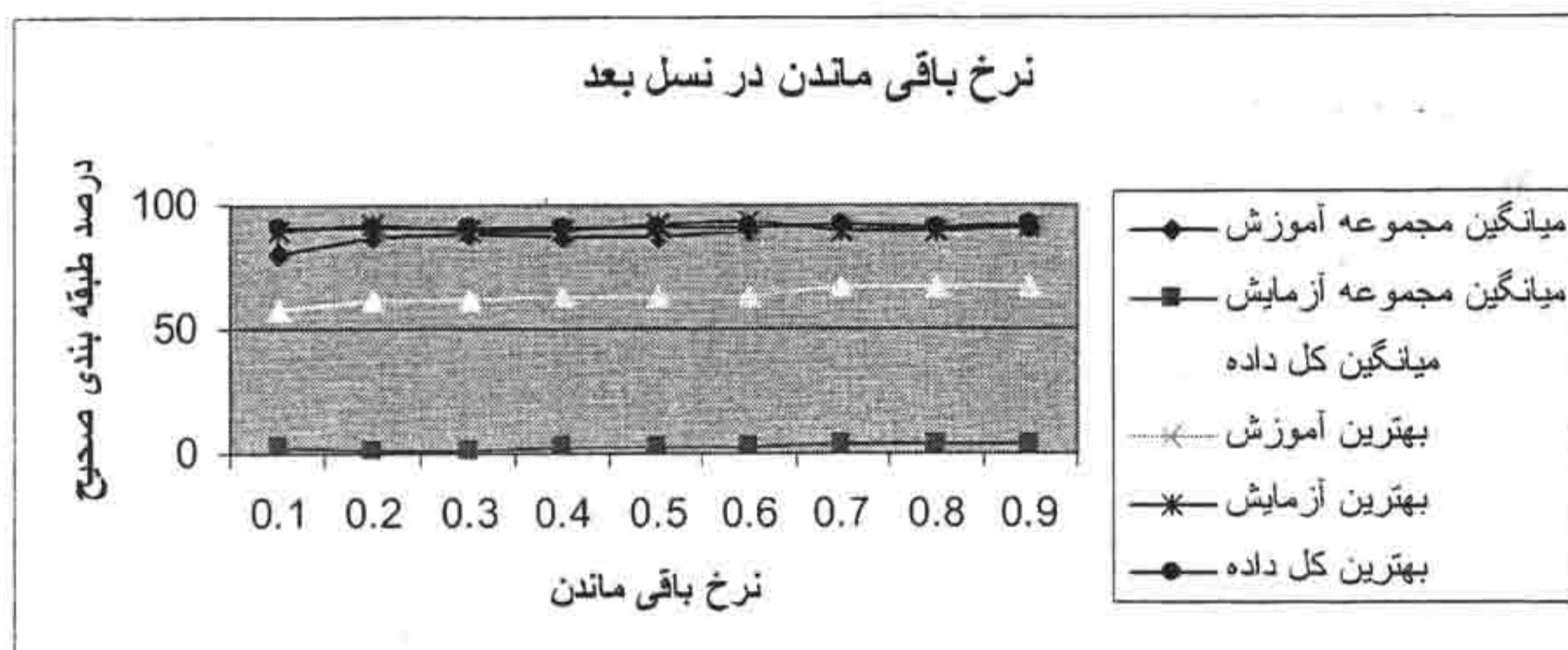
در جدول (۱) تعداد مشخصه ها در حالت اولیه، پس از گسته سازی و پس از انتخاب مشخصه ها نشان داده شده است.

جدول ۱ - تغییرات تعداد مشخصه ها پس از گسته سازی و انتخاب مشخصه های مناسب

درصد کاهش تعداد مشخصه ها	تعداد مشخصه پس از انتخاب	تعداد مشخصه پس از گسته سازی	تعداد مشخصه اولیه	نام مجموعه داده
٪۷۸	۴۴۵	۱۹۹۶	۱۵۵۸	<i>Adeater</i>
٪۷۲	۷۸	۲۸۲	۸۵	<i>Insurance</i>

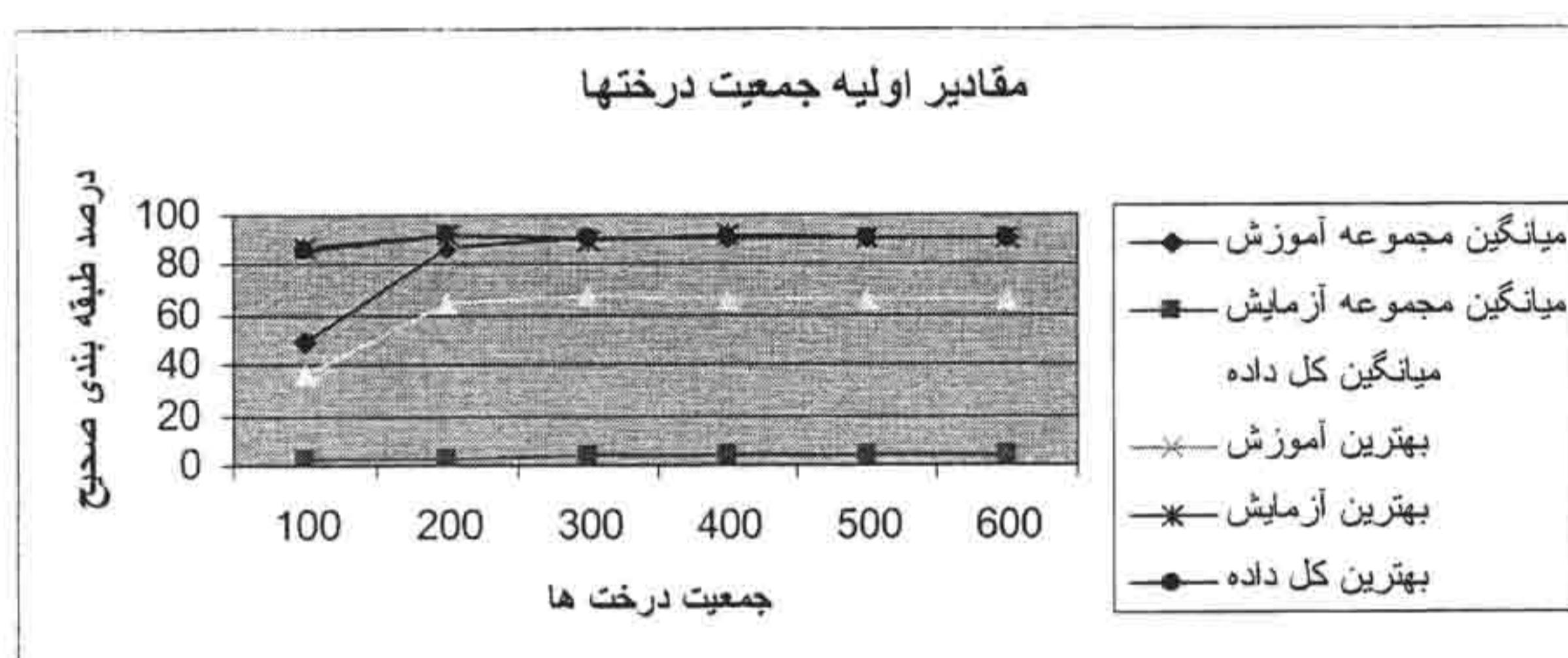
همانطور که مشاهده می شود تعداد مشخصه ها پس از انتخاب مشخصه های مناسب، کاهش چشمگیری داشته است.

نمودارهای مقایسه کارایی پردازش تکاملی بر روی پارامترهای نرخ باقی ماندن در نسل بعد، مقادیر اولیه جمعیت درختها، نقطه قطع و اندازه مجموعه آزمایش بر روی دو مجموعه داده، قبل و بعد از انتخاب مشخصه های مناسب نشان داده شده است.



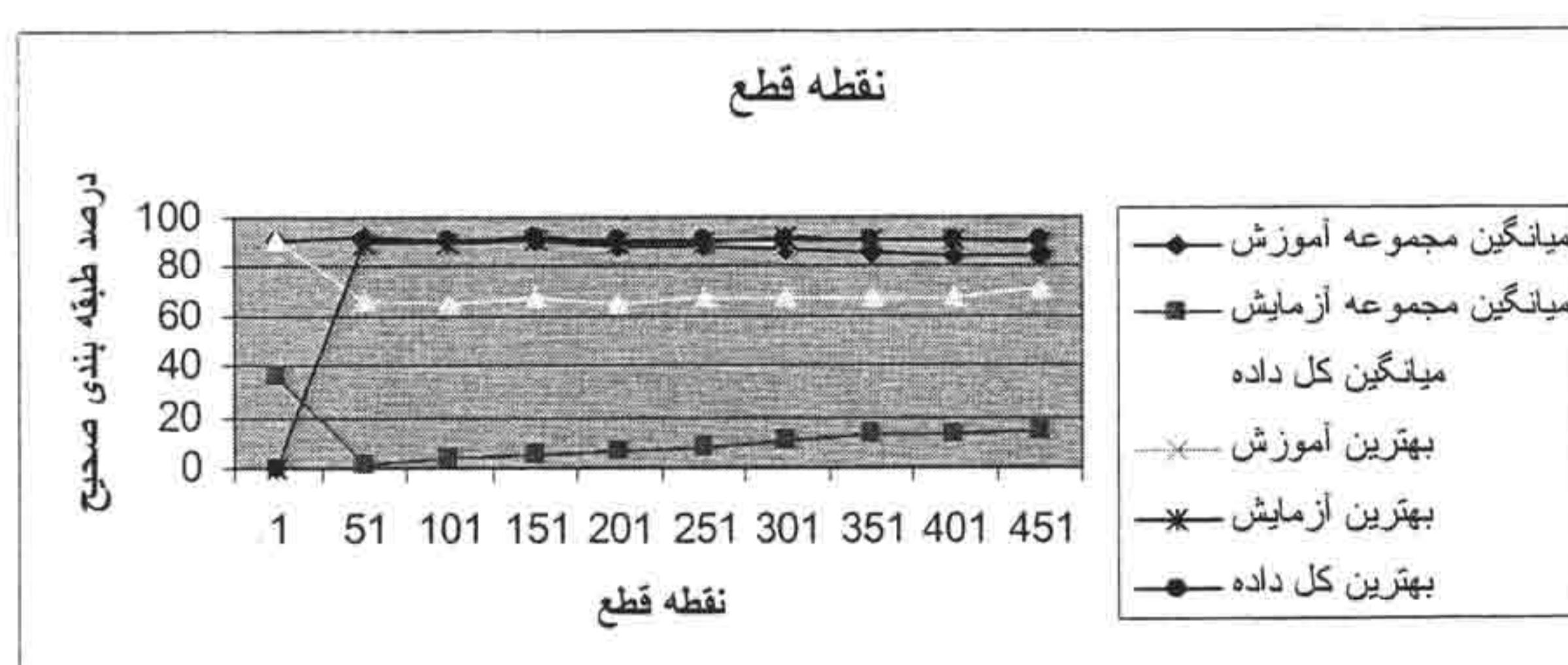
شکل ۲ - نرخ باقی ماندن در نسل بعد برای مجموعه داده Adeater

همانطور که در نمودار فوق مشاهده می شود در نمودار نرخ باقی ماندن در نسل بعد، کارایی با افزایش نرخ باقی ماندن در نسل بعد به آرامی افزایش پیدا میکند.



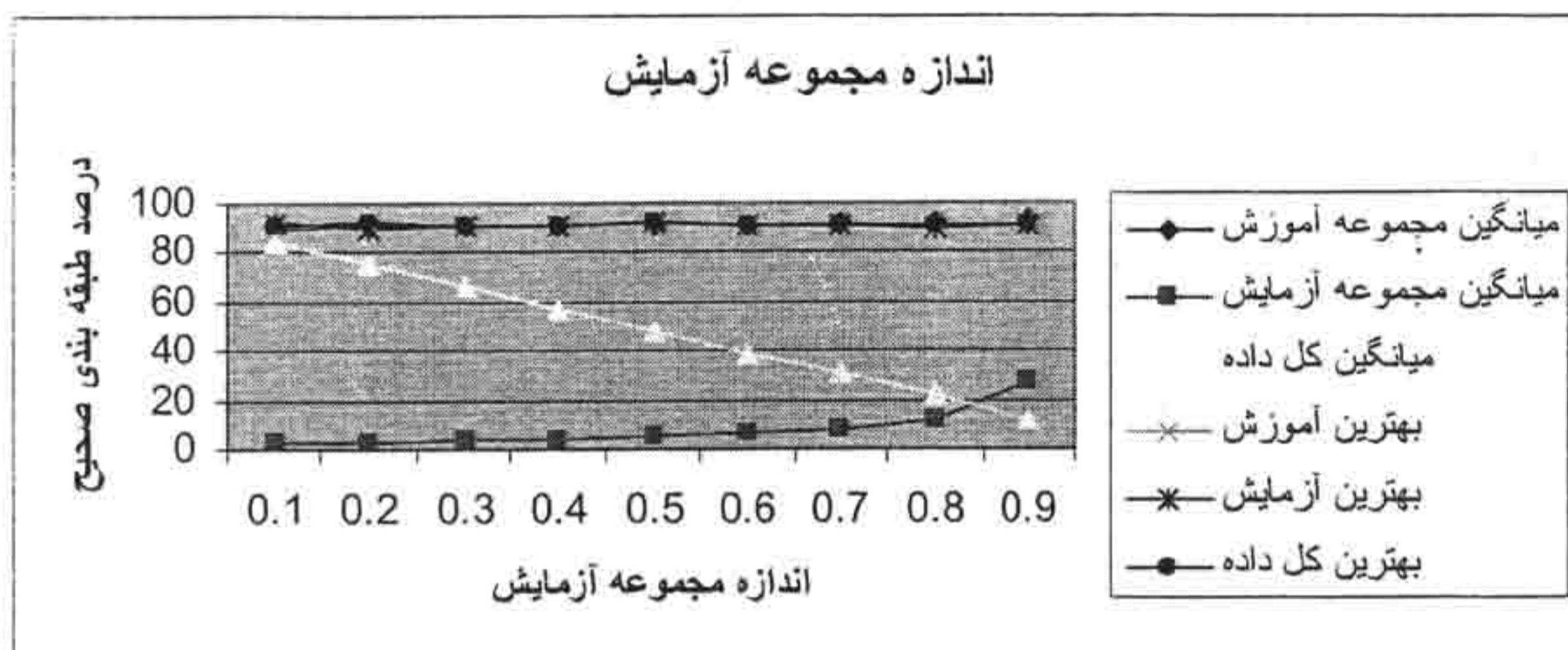
شکل ۳ - مقادیر اولیه جمعیت درخت ها برای مجموعه داده Adeater

در نمودار فوق نتیجه می شود با افزایش جمعیت اولیه درخت ها در صد طبقه بندی افزایش پیدا می کند.



شکل ۴ - نقطه قطع برای مجموعه داده Adeater

نتیجه نمودار فوق این است که با افزایش نقاط قطع در الگوریتم با وجود آنکه از میزان کارایی داده مجموعه آموزشی تا حدود کمی کاسته می شود اما به میزان کارایی مجموعه آزمایش و کل داده بصورت محسوس اضافه می گردد.



شکل ۵ - اندازه مجموعه آزمایش برای مجموعه داده Adeater

از نمودار فوق نتیجه می شود که با افزایش اندازه مجموعه آزمایش ، میزان طبقه بندی صحیح بر روی مجموعه آموزشی کاسته می گردد و با کاهش اندازه مجموعه آزمایش بر میزان طبقه بندی صحیح بر روی مجموعه آموزشی افزوده می گردد و این موضوع منطقی به نظر می رسد زیرا افزایش نمونه های آموزشی باعث تطبیق بیش از حد درخت تصمیم بر روی داده های آموزشی می گردد .

جدول ۲- مقایسه تغییرات طبقه بندی بر روی مجموعه داده های Adeater و Insurance

نام مجموعه داده	تعداد مشخصه اولیه	درصد طبقه بندی حاصل از اجرای الگوریتم C4.5 بر روی تمام مشخصه ها	تعداد مشخصه های انتخابی	درصد طبقه بندی حاصل از اجرای الگوریتم پیشنهادی بر روی مشخصه های انتخابی
Adeater	۱۹۹۶	٪ ۹۷	۴۴۵	٪ ۹۴
Insurance	۲۸۲	٪ ۹۶	۷۸	٪ ۹۴

همانطور که مشاهده می شود با کاهش دادن تعداد مشخصه ها و استفاده از روش پیشنهادی فقط مقدار کمی از درصد طبقه بندی کاسته شده است و این بیانگر این حقیقت است که الگوریتم پیشنهادی بر روی مجموعه داده های بسیار بزرگ بسیار بهینه تر از روش C4.5 عمل خواهد نمود .

۶- نتیجه گیری

در این مقاله یک روش جدید برای تولید درخت تصمیم با استفاده از تئوری مجموعه های دانه درشت و پردازش تکاملی ارائه شد . مراحل مختلف شامل پیش پردازش داده ها ، انتخاب مشخصه های مناسب و ساخت درخت تصمیم به تفصیل بیان شد و نتایج حاصل از این کار با نتایج و الگوریتم های معتبر مورد مقایسه قرار گرفت و نشان داد که این روش در مقایسه با آنها بسیار بهینه عمل می کند . از مزایای این روش می توان به کار ، بر روی مجموعه داده های بزرگ با مجموعه مشخصه های فراوان اشاره نمود و از معایب این روش نیز می توان به کم بودن درصد طبقه بندی آن اشاره نمود . از کارهای آینده برای بهینه نمودن این الگوریتم می توان به جایگزین نمودن روش مورچگان به جای الگوریتم ژنتیک در انتخاب مشخصه ها و ساخت

درخت تصمیم اشاره نمود و از روش‌های بهینه نمودن نتیجه طبقه بندی نظری *Boosting* و *Bagging* برای رسیدن به نتیجه بهتر استفاده نمود.

۷- منابع

- 1) Nicholas kushmerick , *Learning to remove Internet advertisements* ,3rd Int. Conf. on Autonomous Agents 1999.
- 2) Lio,H. And Motoda H., eds., *Instance Selection and Construction for Data Mining* , Kluwer Academic Boston : MA , 2001
- 3) *Coil Challenge 2000* - <http://www.dcs.napier.ac.uk/coil/challenge/>.
- 4) Lian-Yin Zhai*, Li-Pheng Khoo, Sai-Cheong Fok , *Feature extraction using rough set theory and genetic algorithms an application for the simplification of product quality evaluation* , Computers & Industrial Engineering 43 (2002) 661–676
- 5) Komorowski, J., Pawlak, Z., Polkowski, L. & Skowron, A. (1999). *Rough sets: A tutorial*. In *Rough Set Theory*, 3–98.
- 6) D. E. Goldberg. *Genetic algorithms in search, optimization, and machine learning*. Addison-Wesley , Reading, Mass, 1989.
- 7) Jakub Wroblewski, *Finding minimal reducts using genetic algorithm* , Instiute of Mathematics University of Warsaw,1996.
- 8) روش نوین برای بهینه سازی شبکه عصبی با استفاده از پردازش تکاملی. عبدالرضا سالاری ، سعید سهیلی خواه ، پنجمین کنفرانس سیستم های هوشمند ، دانشگاه فردوسی مشهد ۱۳۸۲ (بخش ۲)
- 9) Grefenstette,J.J(1986). *Optimization of control parameters for genetic algorithms*. IEEE Transactions on Systems, Man, and Cybernetics, 16, 122-128.
- 10)M. Kantardzic , *Data Mining Concepts , Models , Method and Algorithms* , Wiley-InterScience , 2003 (Chapter 2)