



دانشگاه شهید بهشتی
دانشکده مهندسی برق و کامپیوتر



ارائه‌ی یک روش جدید برای رده‌بندی مجموعه داده‌های ناهمگون با استفاده از مدل چند عاملی

پایان نامه کارشناسی ارشد مهندسی کامپیوتر
گرایش نرم‌افزار کامپیوتر

استاد راهنما :

دکتر اسلام ناظمی

استاد مشاور:

مهندس حسین علیزاده

ارائه دهنده :

مانی دوستدار

فهرست مطالب

- مقدمه و مفاهیم
- کارهای مرتبط
- روش پیشنهادی
 - ✓ کلیات
 - ✓ سازماندهی قوانین
 - ✓ عملکرد عامل‌ها
- شبیه‌سازی و ارزیابی
- نتیجه‌گیری و کارهای آینده

مقدمه

- امروزه سازمان-ها علاوه بر ذخیره-سازی داده-ها به اطلاعات سطح بالاتری احتیاج دارند تا بتوانند پیش-بینی-هایی انجام داده و تصمیماتی را اتخاذ کنند.
- داده کاوی: داده-کاوی نتیجه-تکامل طبیعی فن-آوری اطلاعات باشد و در واقع به استخراج دانش از حجم بالایی از داده اطلاق می-گردد.

• فعالیت-های داده-کاوی

✓ فعالیتهای توصیفی داده کاوی

✓ فعالیتهای پیشگویانهی داده کاوی

■ رده بندی داده ها



رده بندی

- رده بندی به فرآیند یافتن یک مدل (تابع) گفته می شود که رده های داده ها و مفاهیم را توصیف و مشخص کند.

- مراحل رده بندی داده ها

- ✓ یادگیری

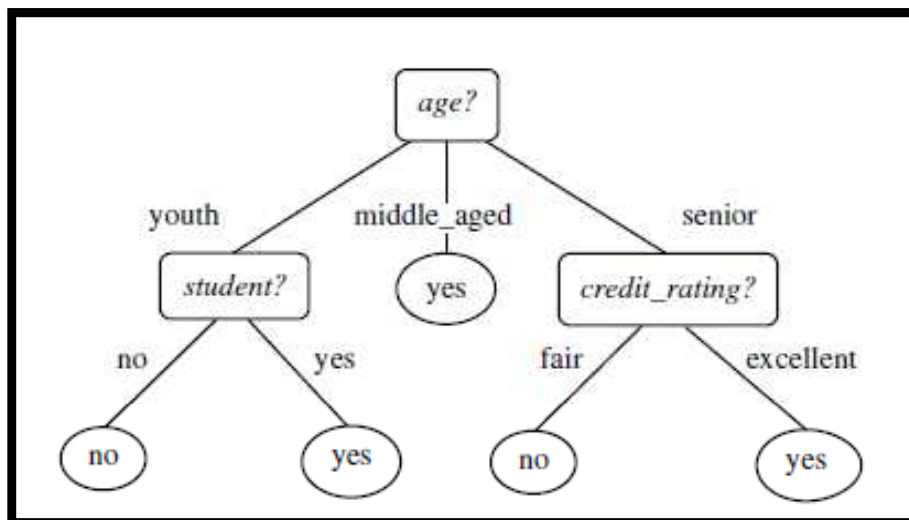
- ✓ آزمایش و استفاده

- نحوه نمایش مدل

- ✓ مجموعه ای از قوانین

- ✓ درخت تصمیم

- ✓ فرمول ریاضی



ارائه ی یک روش جدید برای رده بندی مجموعه داده های ناهمگون با استفاده از مدل چند عاملی

سیستم‌های مبتنی بر عامل

- سیستمی که در طراحی آن از مفهوم عامل استفاده شده است.
- عامل: بخشی از سیستم که به صورت خودکار وظایفی را که بر عهده دارد انجام می‌دهد.
- ویژگی‌ها:
 - ✓ محیط
 - ✓ ورودی و خروجی
 - ✓ سطح اختیار عامل
 - ✓ نحوه ارتباط عامل‌ها
- ناظری بر فعالیت عامل وجود ندارد - داشتن هدف مشترک - ارسال پیام - اهمیت صحت اطلاعات دریافتی
- سطوح مختلف برای عامل‌ها - تلاش برای رسیدن به یک هدف - مدیریت یک سطح از عامل‌ها توسط عامل دیگر

طرح مسئله

- مشکل موجود

✓ در دسترس نبودن مجموعه داده در یک محل و متفاوت بودن مجموعه داده‌ها از نظر برخی از ویژگی‌ها (تعداد صفات و تعداد نمونه‌ها)

- راه حل پیشنهادی

✓ استفاده از سیستم‌های چند عاملی

✓ استفاده از قالب مجموعه قوانین (تشکیل پایگاه جامع قوانین)

ترکیب دو حوزه‌ی داده‌کاوی و سیستم‌های مبتنی بر عامل

- عامل‌های مبتنی بر داده‌کاوی
 - داده‌کاوی مبتنی بر عامل (کاوش مبتنی بر عامل)
- ✓ داده‌کاوی توزیع شده‌ی چند عاملی (استفاده از مدل همتا به همتا)

ارزیابی و استفاده از قوانین

- استفاده از قوانین
 - ✓ استفاده از بهترین قانون
 - ✓ استفاده از K قانون با اولویت بالاتر
 - ✓ استفاده از همه ی قوانین
- مرتب سازی قوانین
 - ✓ روش های مبتنی بر پشتیبانی-صحت (ACS و CSA)
 - ✓ روش های توزین قانون (WRA، LA و χ^2)
- انتخاب قوانین از مجموعه نتایج (حذف قوانین زائد)
 - ✓ قبل از تولید نتایج
 - ✓ بعد از تولید نتایج

ترکیب روش های رده بندی

- در دسترس نبودن تمام حجم داده
✓ ترکیب روش های رده بندی (یادگیری گروهی)
- حجم بالای مجموعه داده
✓ بررسی بخش هایی از مجموعه داده
✓ ترکیب نتایج
- کاهش پیچیدگی و افزایش کارایی (داده با ابعاد بالا)
✓ تقسیم افقی و عمودی

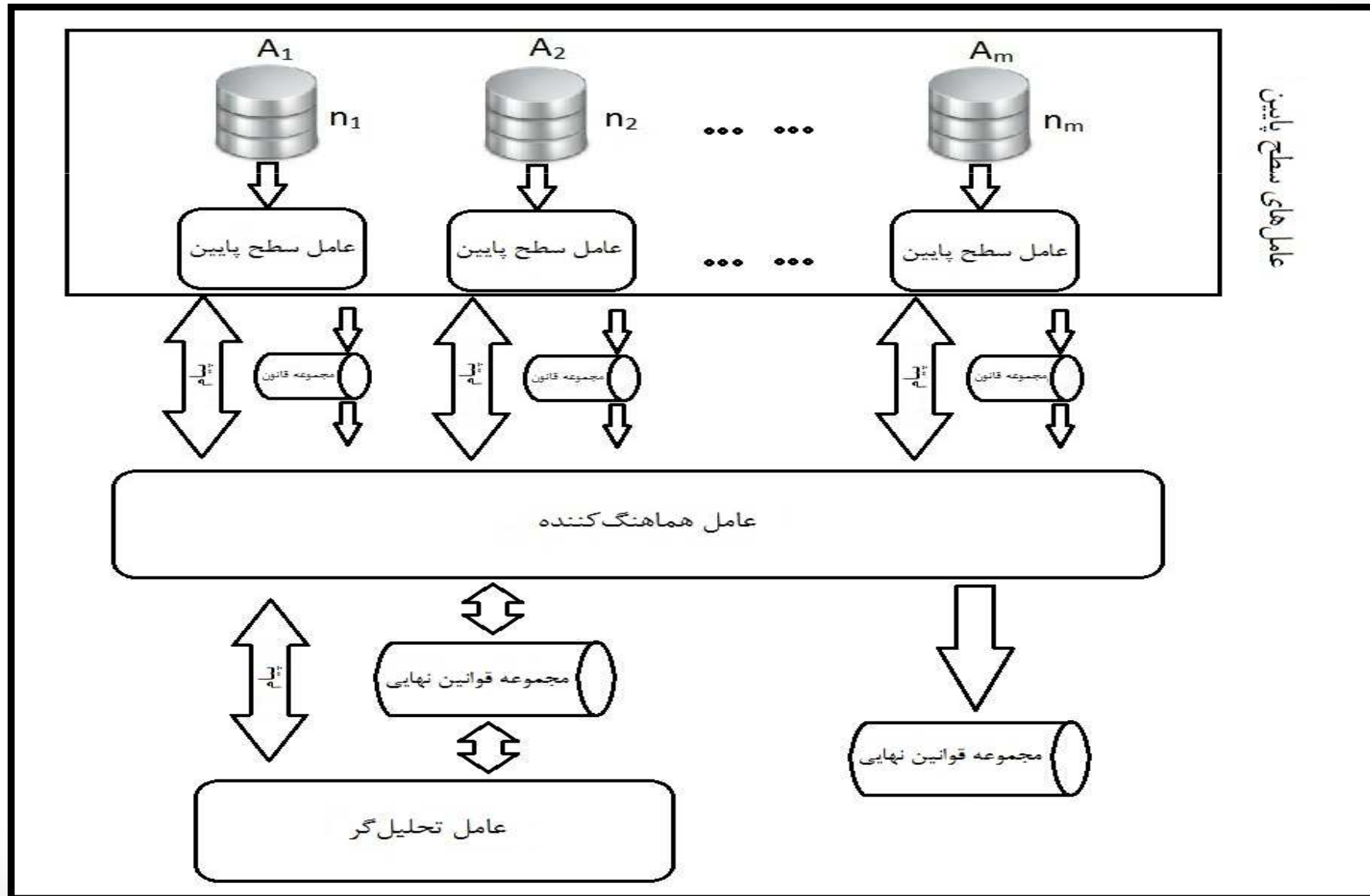
کلیات و ویژگی های اصلی

- ✓ استفاده از سیستم های مبتنی بر عامل برای تحلیل بهتر مسئله (نامتمرکز بودن داده ها و ناهمگونی آنها)
- ✓ بررسی بخشی از داده های حوزه مذکور به منظور تحلیل حوزه و مشخص کردن پارامترهای ارزیابی
- ✓ حضور انواع مختلف عامل ها به منظور انجام وظایف گوناگون در سیستم
- ✓ قالب نتایج به صورت مجموعه قوانین است (نتایج میانی و نهایی)
- ✓ اعتبارسنجی نتایج به صورت دوره ای

عامل ها

- عامل های سطح پایین (Low level)
 - ✓ رده بندی اولیه ی داده ها و ارسال نتایج به عامل هماهنگ کننده
 - ✓ عدم ارتباط با عامل های دیگر در همان سطح
- عامل هماهنگ کننده (Coordinator)
 - ✓ دریافت نتایج از عامل های سطح پایین و تولید مجموعه قوانین نهایی
 - ✓ ارسال قوانین به عامل تحلیل گر برای اعتبارسنجی
- عامل تحلیل گر (Analyst)
 - ✓ دریافت نتایج از عامل هماهنگ کننده و اعتبارسنجی آنها
 - ✓ حذف قوانین زائد

نحوه ارتباط عامل ها با یکدیگر



ارائه ی یک روش جدید برای رده بندی مجموعه داده های ناهمگون با استفاده از مدل چند عاملی

سازماندهی قوانین

- فاکتورهای مهم در رتبه بندی قوانین
- روش محاسبه ی امتیاز قوانین
- تعیین پارامترهای مهم برای ارزیابی قوانین

فاکتورهای مهم در رتبه بندی قوانین

- در سطح عامل سطح پایین
 - ✓ صحت (Confidence)
 - ✓ پشتیبانی (Support)
 - ✓ اندازه ی بخش پیشین (size of Antecedent)
- در سطح عامل هماهنگ کننده
 - ✓ تعداد نمونه ها
 - ✓ تعداد عامل های تولید کننده قانون

روش محاسبه‌ی امتیاز قوانین

- محاسبه‌ی امتیاز جزئی در عامل سطح پایین
✓ صحت قانون (C)

$$C = \frac{n_{true}}{n_{support}}$$

✓ پشتیبانی قانون (S)

$$S = \frac{n_{support}}{n_{total}}$$

✓ نسبت تعداد صفتهای بخش پیشین (A)

$$A = \frac{\text{number of attributes}_{rule}}{\text{number of attributes}_{total}}$$

✓ محاسبه‌ی امتیاز جزئی

$$score_{partial} = w_1 C + w_2 S + w_3 A$$

$$0 < w_i \leq 1$$

روش محاسبه‌ی امتیاز قوانین (ادامه)

- اطلاعات قوانین

- ✓ متن قانون (شرط و بدنه‌ی اصلی)
- ✓ امتیاز جزئی محاسبه شده توسط عامل مربوطه
- ✓ تعداد نمونه‌هایی که قانون از آنها تولید شده است
- ✓ شماره‌ی عامل

روش محاسبه‌ی امتیاز قوانین (ادامه)

- بررسی وضعیت قانون جدید در مقایسه با مجموعه قوانین نهایی
 - ✓ شباهت و تناقض
 - ✓ حذف قوانین زائد
- محاسبه‌ی امتیاز کلی در عامل هماهنگ کننده
 - ✓ نسبت تعداد نمونه‌های عامل

$$0 < \frac{n_{agent}}{n_{max}} \leq 1$$

✓ محاسبه‌ی امتیاز کلی قانون

$$Score_{final} = \frac{\sum_{i=1}^{number\ of\ agents} n_i \times score_i}{number\ of\ agents \times n_{max}}$$

تعیین پارامترهای مهم برای ارزیابی قوانین

- رتبه بندی قوانین
 - ✓ مشخص کردن ضرایب برای محاسبه ی امتیاز قانون توسط هر عامل سطح پایین
- $$score_{partial} = w_1 C + w_2 S + w_3 A$$
- $$0 < w_i \leq 1$$
- انتخاب قوانین
 - ✓ مشخص کردن مقدار آستانه برای فاکتورهای صحت و پشتیبانی به منظور انتخاب قوانین

تعیین پارامترهای مهم برای ارزیابی قوانین (رتبه بندی قوانین)

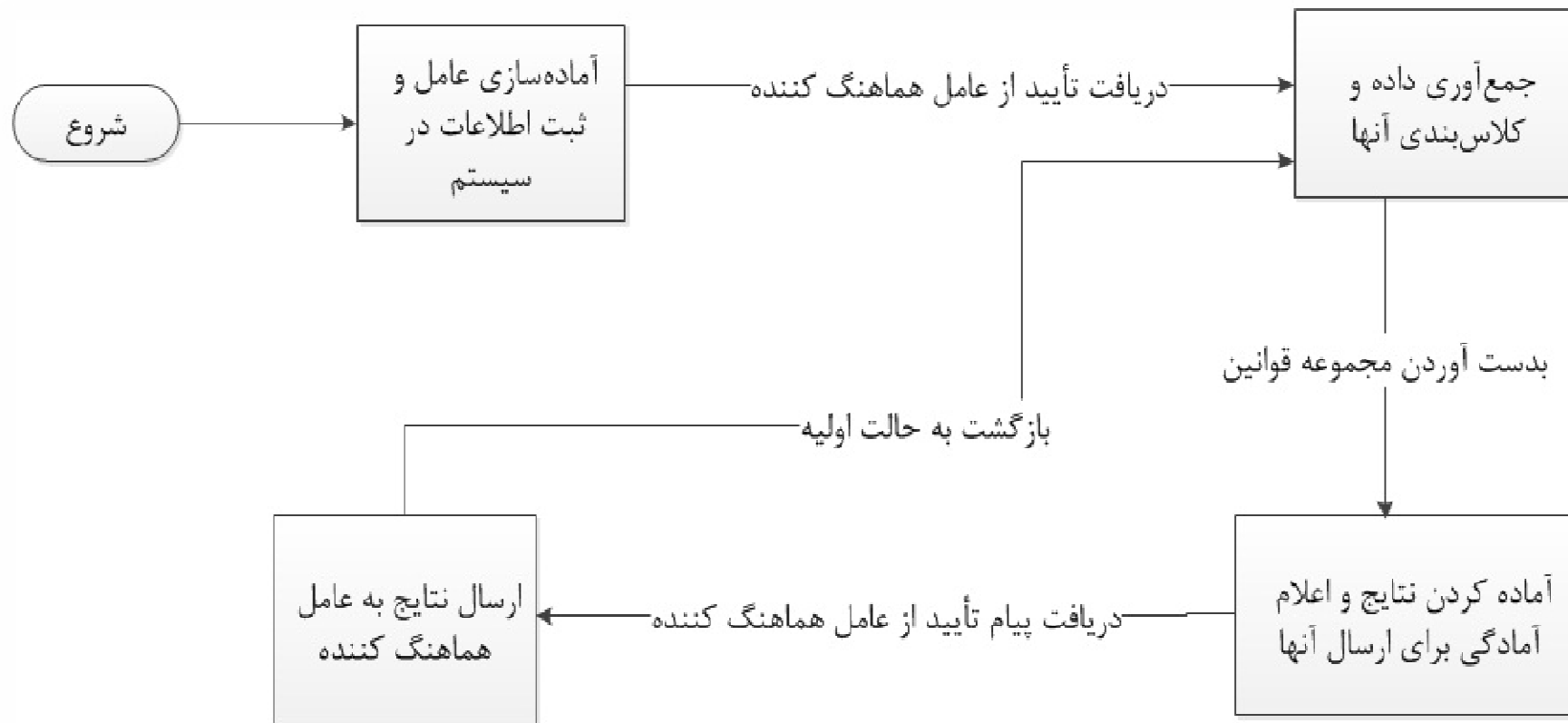
- طراحی آزمایش برای تعیین ضرایب به منظور رتبه بندی قوانین
 - ✓ انتخاب یک مجموعه داده برای تحلیل حوزه (یادگیری و آزمایش)
 - ✓ رده بندی مجموعه داده و بدست آوردن نتایج در قالب مجموعه قوانین
 - ✓ تعیین مقادیر مناسب برای ضرایب رابطه‌ی محاسبه کننده امتیاز قانون توسط عامل سطح پایین (W_i ها)
 - ✓ رتبه بندی نتایج حاصل از رده بندی داده‌ها با استفاده از مقادیر مشخص شده برای ضرایب
 - ✓ آزمایش نتایج رتبه بندی شده با استفاده از مجموعه داده آزمایشی انتخاب شده
 - ✓ انتخاب بهترین ترکیب ضرایب

تعیین پارامترهای مهم برای ارزیابی قوانین (انتخاب قوانین)

- طراحی آزمایش برای تعیین مقادیر آستانه به منظور حذف قوانین زائد
 - ✓ تعیین مقادیر مناسب برای آستانه دو فاکتور صحت و پشتیبانی
 - ✓ انتخاب یک ترکیب از مقادیر آستانه و آزمایش مجموعه قوانین انتخاب شده با استفاده از مجموعه داده‌های آزمایشی در هر تکرار
 - ✓ انتخاب بهترین ترکیب برای مقادیر آستانه

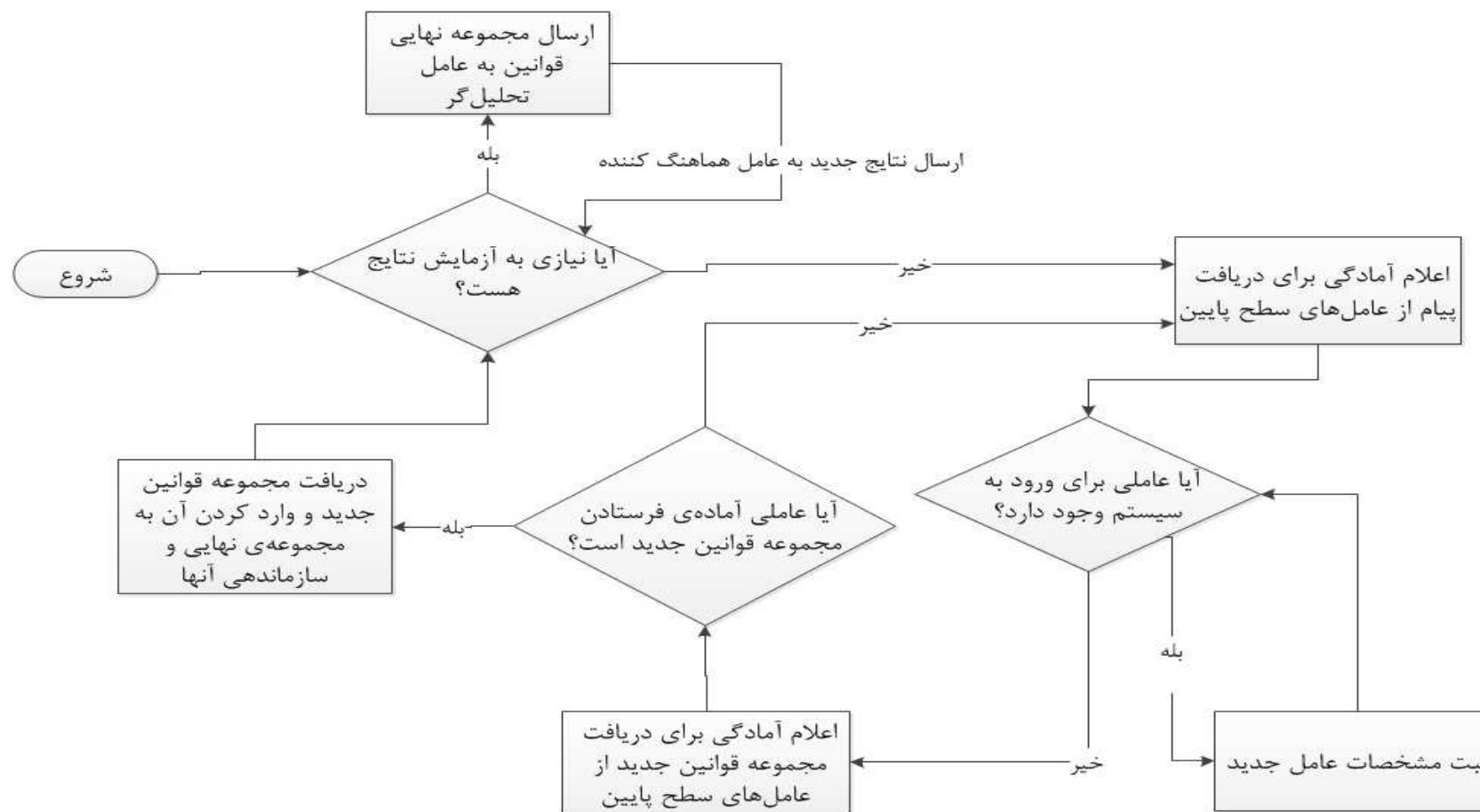
شرح عملکرد عامل ها

• عامل های سطح پایین



شرح عملکرد عامل‌ها

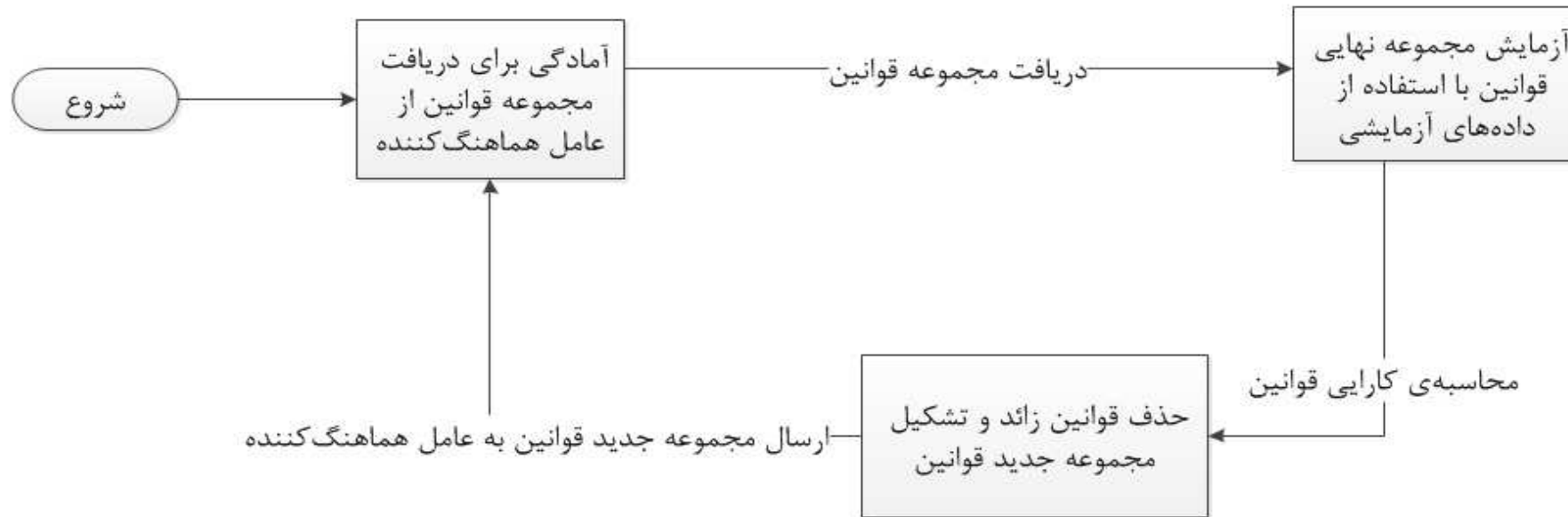
• عامل هماهنگ کننده



ارائه‌ی یک روش جدید برای رده‌بندی مجموعه داده‌های ناهمگون با استفاده از مدل چند عاملی

شرح عملکرد عامل‌ها

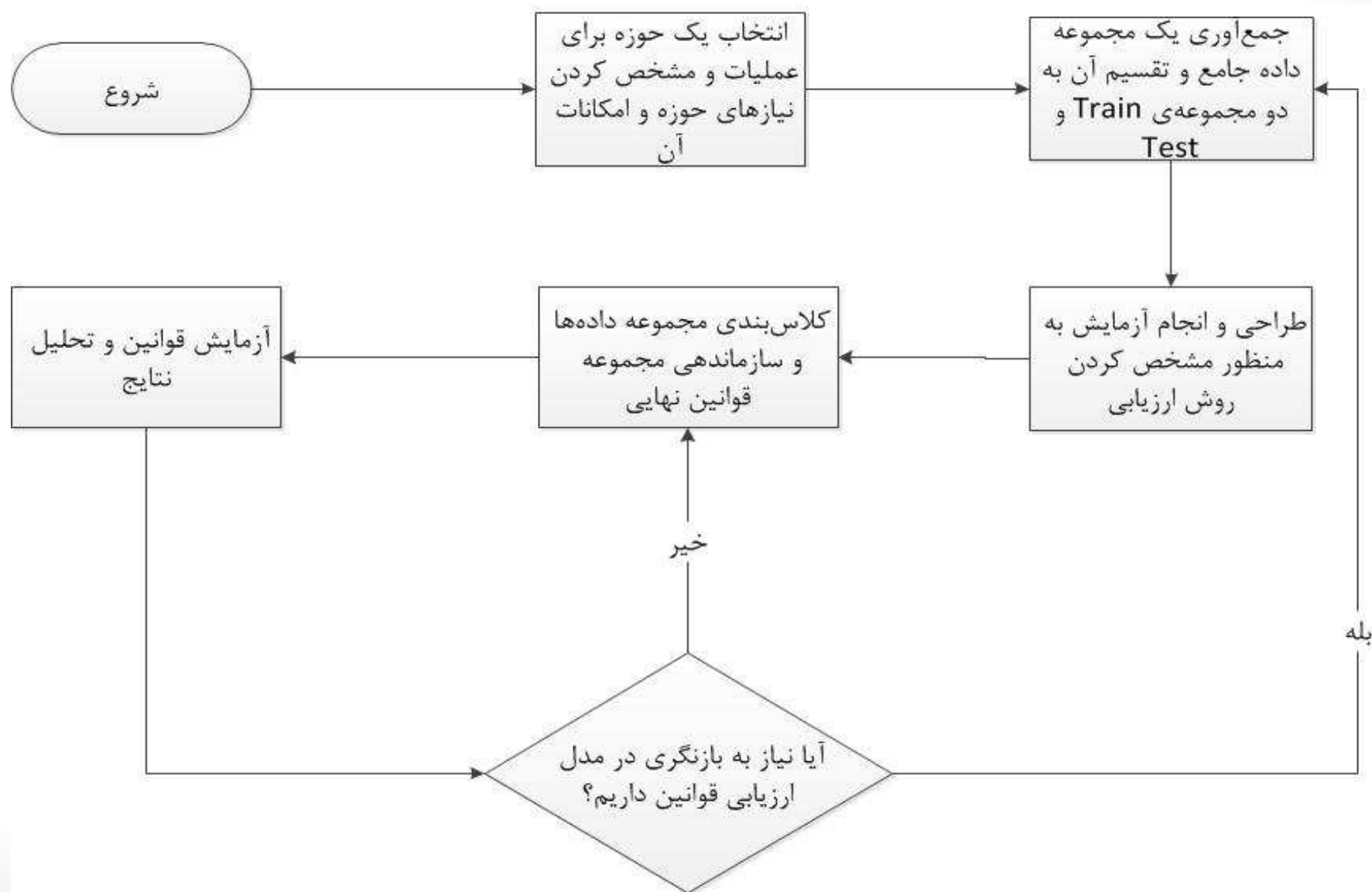
• عامل تحلیل‌گر



$$Efficiency = \frac{n_{correctly\ classified} - n_{incorrectly\ classified}}{n_{classified}}$$

$$-1 \leq Efficiency \leq 1$$

نمای کلی روش پیشنهادی



ارائه ی یک روش جدید برای رده بندی مجموعه داده های ناهمگون با استفاده از مدل چند عاملی

ابزارهای مورد استفاده



IKVM.NET

- مجموعه نرم‌افزاری WEKA
 - ✓ مجموعه‌ای از الگوریتم‌های یادگیری ماشین
 - ✓ ابزارهای پیش پردازش داده‌ها
 - ✓ نوشته شده به زبان جاوا

- بسته نرم‌افزاری Microsoft visual studio
 - ✓ پیاده‌سازی عامل‌ها

✓ استفاده از زبان برنامه‌سازی C#.net

- بسته‌ی نرم‌افزاری IKVM

✓ رابط بین عامل‌ها و نرم‌افزار WEKA

✓ تبدیل کلاس‌های جاوا به فایل‌ها DLL

شبیه سازی عامل ها

- برقراری ارتباط
 - ✓ عامل هماهنگ کننده دخالتی در وظایف عامل های سطح پایین ندارد
 - ✓ عامل ها تنها با ارسال پیام با یکدیگر ارتباط برقرار می کنند
- عامل سطح پایین
 - ✓ استفاده از روش PART برای رده بندی اولیه ی داده ها (دلیل: قالب نتایج الگوریتم)
 - ✓ تبدیل نتایج و اطلاعات تولید شده به یک قالب خاص برای پردازش آسان تر

روش ارزیابی

- مقایسه‌ی روش پیشنهادی با روش‌هایی که مجموعه داده را به صورت یکجا بررسی می‌کنند
- استفاده از یک مجموعه داده و تقسیم آن برای مراحل مختلف روش پیشنهادی ✓
 - مجموعه داده برای تحلیل حوزه (یادگیری و آزمایش)
 - مجموعه داده برای انجام عملیات رده‌بندی و سازماندهی مجموعه نهایی قوانین (یادگیری و آزمایش)
 - ❖ تقسیم داده‌ها بین عامل‌های مختلف در تکرارهای متفاوت
- ✓ روش‌های مورد مقایسه
 - مجموعه داده برای انجام عملیات رده‌بندی (یادگیری و آزمایش)

مجموعه داده‌ی Adult

- مشخصات

تعداد صفت‌ها	نوع صفت‌ها	تعداد نمونه برای یادگیری	تعداد نمونه برای آزمایش	تعداد نمونه برای تحلیل حوزه
۱۲	مقدار قطعی	۳۲۵۶۰	۱۱۲۸۲	۵۰۰۰

- نتایج تحلیل حوزه

روش	CSA	ACS	New method
درستی روش	٪۷۶	٪۷۳	٪۸۰

✓ آستانه‌ی صحت: ۰.۷

✓ آستانه‌ی پشتیبانی: ۰.۰۲

$$score_{partial} = 0.2 \times C + S + 0.3 \times A$$

مجموعه داده‌ی Adult

- درستی روش در پایان تکرارها

تکرار سوم	تکرار دوم	تکرار اول	روش / تکرار
۸۳.۱۹	۸۳.۲۹	۸۱.۴۰	درستی روش

- مقایسه‌ی نتایج

روش	روش	روش	روش	روش	روش	روش	روش	روش	روش
J48	JRip	Decision Table	Ridor	CR	ZeroR	Random Forest	PART	پیشنهادی	روش
۸۱.۹۱	۸۲.۵۳	۸۱.۷۸	۸۱.۶۵	۷۵.۶۴	۷۵.۶۴	۸۱.۹۷	۸۳.۱۴	۸۳.۱۹	درستی

مجموعه داده‌ی Car

- مشخصات

تعداد صفت‌ها	نوع صفت‌ها	تعداد نمونه برای یادگیری	تعداد نمونه برای آزمایش	تعداد نمونه برای تحلیل حوزه
۷	مقدار قطعی	۹۶۳	۴۵۹	۳۰۶

- نتایج تحلیل حوزه

روش	CSA	ACS	New method
درستی روش	%۸۳.۰۱	%۸۳.۸۵	%۸۳.۶۶

$$score_{partial} = 0.2 \times C + S + 0.5 \times A$$

مجموعه داده‌ی Car

- درستی روش در پایان تکرارها

تکرار اول	تکرار دوم	روش / تکرار
۷۸	۸۵.۰۸	صحت روش

- مقایسه‌ی نتایج

روش پیشنهادی	روش PART	روش Random Forest	روش ZeroR	روش CR	روش Ridor	روش Decision Table	روش JRip	روش J48	روش
۸۵.۰۸	۸۵.۷۱	۶۰.۶۱	۶۷.۵۶	۶۷.۹۳	۷۹.۳۰	۷۶.۰۶	۶۴.۸۶	۷۶.۴۴	درستی

Nursery داده‌ی

- مشخصات

تعداد صفت‌ها	نوع صفت‌ها	تعداد نمونه برای یادگیری	تعداد نمونه برای آزمایش	تعداد نمونه برای تحلیل حوزه
۹	مقدار قطعی	۷۴۶۰	۳۵۰۰	۲۰۰۰

- نتایج تحلیل حوزه

روش	CSA	ACS	New method
درستی روش	%۹۲.۴	%۷۶.۸	%۹۶

$$score_{partial} = 0.2 \times C + S + 0.5 \times A$$

Nursery داده‌ی

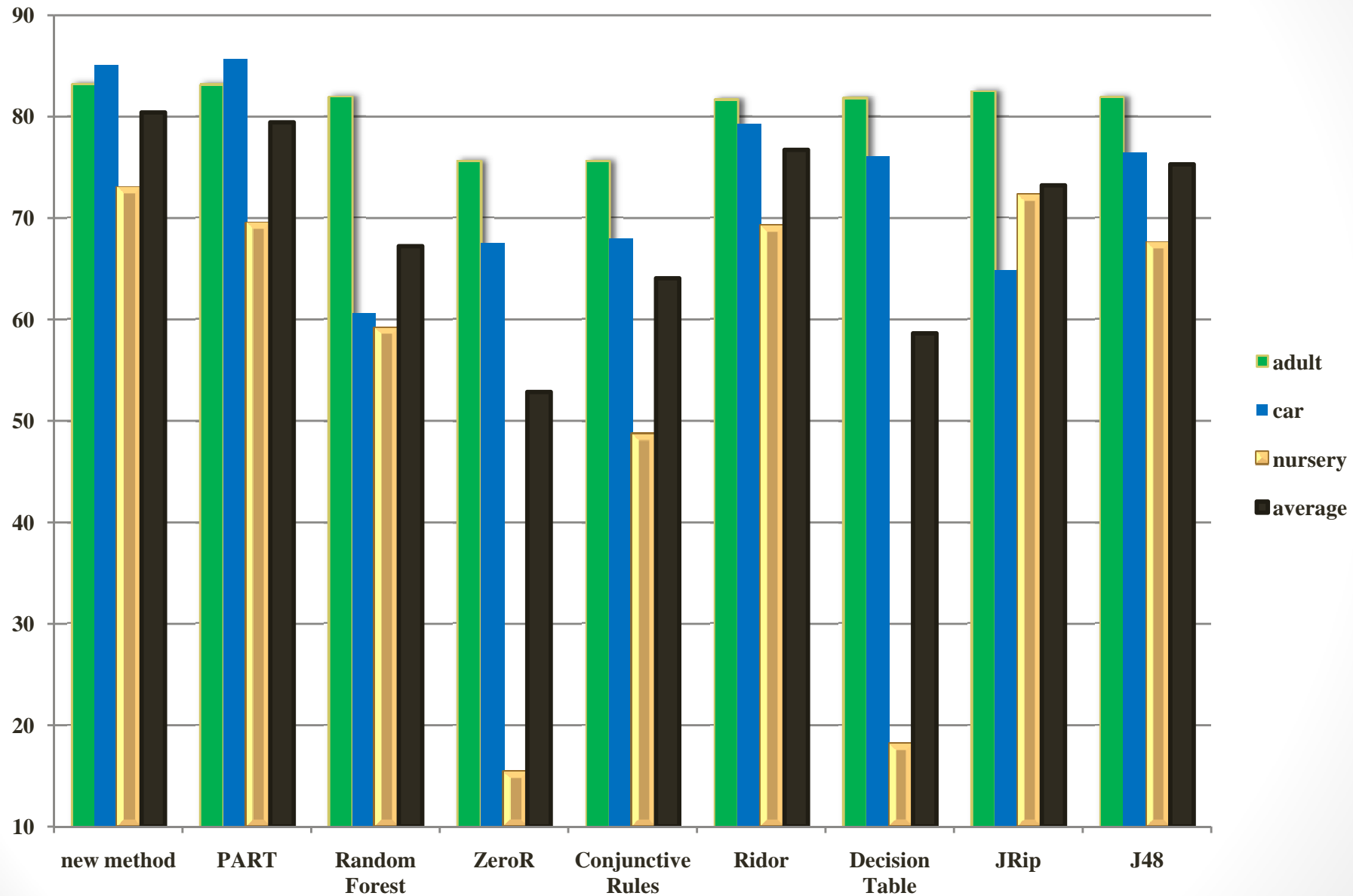
- درستی روش در پایان تکرارها

تکرار دوم	تکرار اول	روش / تکرار
۷۲.۹۵	۵۹.۲	صحت روش

- مقایسه‌ی نتایج

روش	روش	روش	روش	روش	روش	روش	روش	روش	روش
J48	JRip	Decision table	Ridor	CR	ZeroR	Random Forest	PART	پیشنهادی	روش
۶۷.۵۵	۷۲.۳	۱۸.۱	۶۹.۲	۴۸.۶۵	۱۵.۳۵	۵۹.۱	۶۹.۵	۷۲.۹۵	درستی

مقایسه‌ی نتایج



ارائه‌ی یک روش جدید برای رده‌بندی مجموعه داده‌های ناهمگون با استفاده از مدل چند عاملی

تحلیل و بررسی نتایج

- کارایی مطلوب در مورد دو مجموعه داده‌ی Nursery و Adult
 - ✓ استفاده از نتایج همه‌ی عامل‌ها و سازماندهی آنها در قالب یک مجموعه‌ی نهایی
 - ✓ بکارگیری یک روش مناسب برای رتبه‌بندی قوانین (با استفاده از تحلیل حوزه)
 - ✓ حذف قوانین زائد در پایان هر تکرار
- کارایی نامطلوب در مورد مجموعه داده‌ی Car
 - ✓ وابستگی روش مربوطه به مجموعه داده مورد بررسی
 - کم بودن حجم داده به منظور تحلیل حوزه و آزمایش نتایج
 - نامتوازن بودن مجموعه داده‌ی مورد بررسی

نتیجه گیری

• برتری ها

- ✓ نبود مشکلات انتقال و ذخیره سازی داده
- ✓ بکارگیری روش مربوطه بدون ایجاد اختلال در عملکرد سیستم موجود
- ✓ انعطاف پذیری بیشتر در رتبه بندی و انتخاب قوانین (تحلیل حوزه با استفاده از یک مجموعه داده در همان حوزه)
- ✓ در اختیار داشتن یک مجموعه ی بهینه از قوانین (رتبه بندی مناسب و انتخاب قوانین مفید)

• کاستی ها

- ✓ وابستگی به مجموعه داده ی منبع

کارهای آینده

- استفاده از فاکتورهای بیشتر به منظور ارزیابی قوانین
 - ارزیابی و انتخاب قوانین در عامل‌های سطح پایین
 - پوشش انواع متفاوت از داده‌ها
- ✓ بررسی داده‌ها به منظور کشف دنباله‌ای از وقایع (کاربرد در حوزه‌ی پزشکی)

مقالات مستخرج از پایان نامه

- دوستدار، م.، ناظمی، ا. و علیزاده، ح.، ”ارائه‌ی یک روش جدید برای رده‌بندی مجموعه داده‌های ناهمگون با استفاده از مدل چند عاملی“، کنفرانس فناوری اطلاعات و دانش (ارسال شده).

منابع منتخب

- Han, J. and Kamber, M., Data mining: concepts and techniques, Elsevier Pub., 2006.
- Cao, L., Data mining and multi-agent integration, Springer Science+Bussines Media, 2009.
- Gorodetsky, V., Autonomus intelligent systems: agents and data mining, Springer-Verlag, 2007.
- Rao, K. R., Rao, M. N., Prasad Y. S. and Ramya, K. R., “Interaction and Integration of Agent Mining in Distributed Data Environment”, In *International Journal of Computer Science & Emerging Technologies* (E-ISSN: 2044-6004), Volume 1, Issue 4, 2010.
- Gao, J., Denzinger, J. and James, R. C., “A Cooperative Multiagent Data Mining Model and Its Application to Medical Data on Diabetes”, In *Proc. Autonomous Intelligent Systems: Agents and Data Mining*, LNAI 3505, pp. 93–107, Springer-Verlag, 2005.
- Veloso, A. and Meira Jr, W., “Rule generation and rule selection techniques for cost-sensitive associative classification.” *Proc. of the Brazilian Symposium on Databases (SBDD)*. 2005.
- Wang, Y. J., Xin, Q. and Coenen, F., “Hybrid rule ordering in classification association rule mining.” *Transactions on Machine Learning and Data Mining* 1.1, pp. 1-16, 2008.
- Coenen, F. and Leng, P., “An evaluation of approaches to classification rule selection.” *Data Mining, 2004. ICDM'04. Fourth IEEE International Conference on*. IEEE, 2004.
- Baralis, E., Chiusano, S. and Garza, P., “On support thresholds in associative classification”, In *Proceedings of the 2004 ACM symposium on Applied computing*, pp. 553-558. ACM, 2004.
- Todorovski, L. and Džeroski S., “Combining classifiers with meta decision trees.” *Machine Learning* 50, pp. 223-249, 2003.

با سپاس فراوان از همراهی و توجه شما

