

کاربرد تئوری مجموعه های راف در نظریه تصمیم گیری

آصف زارع

استادیار گروه مهندسی برق

دانشگاه آزاد اسلامی واحد گناباد

غلامرضا نساجیان

دانشجوی کارشناسی ارشد مهندسی کنترل

دانشگاه آزاد اسلامی واحد گناباد

چکیده: در این مقاله تئوری مجموعه های راف که تعمیمی از نظریه های مجموعه کلاسیک بر پایه منطق سه مقداری، جهت کار باداده های ناقص و ناسازگار و کاهش داده های مازاد بر نیاز پایگاه داده ها است معرفی می شود. پس از ارائه تاریخچه مختصری از نحوه ارائه شدن این نظریه، مفهوم مجموعه راف و مزیت های استفاده از آن ذکر خواهد شد. آنگاه به طور مختصر تئوری مجموعه راف در آنالیز منطقی داده و استخراج حداقل مجموعه قواعد "اگر... آنگاه..." ارائه شده و در پایان کاربرد تئوری مجموعه راف در آنالیز داده که شامل کشف وابستگی ها در داده و حذف قسمتهای اضافی از داده است با مثال واقعی ارائه می شود.

واژه های کلیدی: مجموعه های راف، سیستم های اطلاعاتی، کاهش داده ها، ویژگیهای تصمیم، ویژگیهای موقعیت

مقدمه:

یافتن یک واژه معادل برای عبارت "ROUGH SETS" کار مشکلی است. در فرهنگ لغت برای کلمه "ROUGH" معادلهایی مانند زبر، درشت، تقریبی، بی ادب، متلاطم و ناصاف در نظر گرفته شده است [1] که در میان این واژه ها کلمه تقریبی شباهت بیشتری با مفهوم مورد نظر بنیانگذار این نظریه دارد. اما هیچ یک از این کلمات بار معنایی خود واژه لاتین را ندارند به همین علت در این مقاله عبارت "مجموعه های راف" به عنوان معادل آن به کار گرفته خواهد شد. تئوری مجموعه های راف در اوایل سال ۱۹۸۰ میلادی توسط پروفیسور زدیسلو پاولاک پایه گذاری شد. این تئوری با تحلیل جدولهای داده سروکار دارد. در این تئوری جدولهای داده می توانند توسط اندازه گیری یا افراد متخصص و آگاه (خبره) بدست آمده باشد. هدف اصلی از تحلیل مجموعه راف به دست آوردن مفاهیم تقریبی از داده های اکتسابی می باشد. این تئوری، یک ابزار قدرتمند ریاضی برای استدلال در موارد ابهام و نایقینی است که روشهایی را برای زدودن و کاستن اطلاعات دانش نامربوط یا مازاد بر نیاز از پایگاههای داده ها مهیا می سازد. این فرایند حذف داده های زائد، بر مبنای آموزش وظیفه اصلی سیستم، وبدون از دست دادن داده های اساسی پایگاه داده ها صورت می پذیرد. در نتیجه تقلیل اطلاعات، مجموعه ای از قواعد تلخیص شده و پرمعنا حاصل می گردد که کار تصمیم گیرنده را بسیار ساده تر می کند. در حقیقت میتوان گفت که مجموعه راف با کاهش فضای داده ها و بر گزیدن ترمهای مهم، یک نگاهت از فضای داده های خام و ترمها به فضای سمانتیک (مفاهیم) انجام می دهد. لذا، با توجه به رشد انفجاری حجم اطلاعات، مجموعه راف می تواند نقش بسیار موثری در سیستمهای پشتیبان تصمیم گیری

داشته باشند [2]. تئوری مجموعه راف نقاط اشتراک زیادی با تئوری مجموعه های فازی، تئوری شهود، روش های استدلال بولی و تحلیل تفکیکی دارد. اما تئوری مجموعه راف به عنوان یک تئوری مستقل در نظر گرفته می شود.

سیستمهای اطلاعاتی و مجموعه های راف:

در بیشتر موارد، اطلاعات بصورت جدول های داده ها که به نام سیستمهای اطلاعاتی، جدولهای ویژگی-مقدار و یا جدولهای اطلاعات در اختیار است. ستون های جدول اطلاعاتی با ویژگیها، وسطهای آن با شیء ها نام گذاری می شوند و ورودی های جدول با مقادیر آن ویژگیها پر می شوند. شیء های دارای مقادیر ویژگیهای یکسان - از نظر آن ویژگی - مشابه تلقی می شوند و به بلوک یکسانی از تقسیم بندی (کلاسه سازی) - که توسط آن مجموعه از ویژگی ها تعیین می شوند - تعلق می یابند. میتوان از مجموعه راف در حل مسائل اساسی در زمینه تحلیل داده ها استفاده نمود، از جمله [3,4]:

* مشخص کردن مجموعه ای از اشیاء بر حسب مقادیر ویژگیها

* یافتن وابستگیها بین ویژگیها

* زدودن (کاهش یا تقلیل) ویژگیهای مازاد (داده ها)

* یافتن مهمترین ویژگیها

* تولید قواعد تصمیم گیری

تئوری مجموعه راف یک روش ریاضی جدید برای تحلیل داده های هوشمند و داده کاوی می باشد. بعد از تقریباً بیست و اندی سال از پایه گذاری تئوری مجموعه راف روشهای کاربردی از آن به یک درجه مشخص از کمال رسیده است و در چند سال اخیر یک رشد سریع از توجه به تئوری مجموعه راف و کاربردهایش در سرتاسر جهان به وضوح دیده میشود. بسیاری از کارگاههای آموزشی، کنفرانسها و سمینارها در برنامه هایشان مجموعه های راف را در نظر می گیرند. نزدیک به دو هزار مقاله و چندین کتاب تا هم اکنون روی جنبه های مختلف از مجموعه های راف منتشر شده است. فلسفه مجموعه های راف بر این فرض است که هر شیء از جهان را می توان به عنوان اطلاعات (داده، معرفت) در نظر گرفت بنا شده است. اشیاء توصیف شده به وسیله اطلاعات یکسان از نقطه نظر اطلاعات در دسترس درباره آنها غیر قابل تشخیص هستند. رابطه غیر قابل تشخیص بودن (رابطه علی- معلولی) به دست آمده در این روش اساس ریاضیات تئوری مجموعه های راف می باشد. هر مجموعه ای از اشیاء غیر قابل تشخیص را یک مجموعه بنیادی می نامند و شکل یک جزء اصلی (اتم) از دانش درباره جهان است. به هر اجتماعی از مجموعه های بنیادی عنوان مجموعه های کریسپ (دقیق) نسبت می دهند و در غیر این صورت مجموعه مبهم و غیر صریح است که عنوان مجموعه راف برای آن در نظر گرفته می شود. تئوری مجموعه های راف مبتنی بر مفهوم کلاسه سازی (دسته بندی) است. برای روشنتر شدن موضوع، به عنوان مثال، گروهی از بیماران که از یک بیماری معین رنج می برند را در نظر بگیرید. با هر بیمار، یک فایل داده ها شامل اطلاعاتی از قبیل نام، آدرس، سن، جنسیت، دمای بدن، فشار خون و مانند آن - همراه است. تمام بیمارانی که علائم مشابهی را نشان می دهند - از نظر اطلاعات در دسترس - مانند یکدیگر هستند، و می توان آنان را در دسته هایی - به عنوان اجزای بنیادی دانش و معرفت موجود نسبت به بیماران - کلاسه بندی نمود. این اجزاء به نام مجموعه های بنیادی یا مفاهیم بنیادی شناخته میشوند و میتوانند به عنوان بلوک های سازنده دانش در مورد بیماران در نظر گرفته شوند. مفاهیم بنیادی می توانند با مفاهیم مختلط - یعنی مفاهیمی که به شکل یکتایی بر حسب مفاهیم بنیادی تعریف می شوند - ترکیب گردند. هر اجتماعی از مفاهیم بنیادی، مجموعه کریسپ - به معنای مجموعه ای با تعریف و مرزهای دقیق - نامیده می شود، و هر مجموعه دیگری که کریسپ نباشد، مجموعه راف - به معنای مجموعه مبهم و نادقیق شناخته می شود. با هر مجموعه X دو مجموعه کریسپ به نامهای تقریب بالایی و تقریب پایینی مجموعه X مرتبط و همراه می شود. تقریب پایینی X اجتماع تمام مجموعه های بنیادینی است که X شامل آنها می شود، در حالیکه تقریب بالایی X اجتماع تمام عناصری است که مطمئناً به X تعلق دارند، در حالیکه تقریب بالایی مجموعه X ، مجموعه تمام عناصری است که احتمالاً به X تعلق دارند. تفاوت تقریب های بالایی و پایینی مجموعه X در نواحی مرزی آنان می باشد. یک مجموعه راف است اگر ناحیه مرزی نانهی داشته باشد، و در غیر این صورت، کریسپ است. عناصر ناحیه مرزی با اطلاعات در دسترس از خود مجموعه یا مکمل

آن قابل کلاسه سازی نیستند. تقریب مجموعه ها عملیات اصلی تئوری راف ست است که برای تعریف کردن وابستگی ها (کلی یا جزئی) بین ویژگیها، تقلیل (زدودن) ویژگیها تولید قواعد تصمیم و غیره به کار گرفته می شود [3,4,5].

یک مثال ساده برای سیستم های اطلاعاتی در جدول زیر آمده است [6]:

store	E	Q	L	P
1	High	Good	no	profit
2	Med	good	no	less
3	Med	good	no	profit
4	No	ave	no	less
5	Med	ave	yes	less
6	High	ave	yes	profit

در جدول زیر شش فروشگاه بر حسب چهار ویژگی E, Q, L و P توصیف شده است:

E: اختیارات پرسنل فروشگاه Q: کیفیت کالاها و اجناس
L: موقعیت رفت و آمد زیاد P: سود یا ضرر فروشگاه

هر فروشگاه بر حسب چهار ویژگی E, Q, L و P دارای توصیف مختلفی می باشد. بنابراین فروشگاهها را می توان با بکارگیری اطلاعات فراهم شده توسط ویژگیها از یکدیگر متمایز نمود. فروشگاه ۲ و ۳ بر حسب ویژگیهای E, Q و L غیر قابل تشخیص هستند زیرا دارای ارزش مشابه در این ویژگیها می باشند. به همین نحو فروشگاههای ۱ و ۲ و ۳ بادر نظر گرفتن ویژگیهای L و Q غیر قابل تشخیص هستند زیرا دارای ارزش مشابه در این ویژگیها می باشند. می توان فروشگاهها را بر حسب ویژگیهای مشابهی که دارند در یک کلاس جمع کردو به عنوان یک زیر مجموعه در نظر گرفت. به عنوان مثال می توان با توجه به ویژگیهای L و Q همه فروشگاهها را در کلاسهای (زیر مجموعه های) زیر طبقه بندی نمود:

{1,2,3} Q:good , L:no

{4} Q:ave , L:no

{5,6} Q:ave L: yes

سوال: آیا می توان مجموعه {1,3,6} را بر حسب ویژگیهای E, Q و L توصیف کرد؟ بدیهی است که نمی توان به صورت منحصر به فرد به این سوال جواب داد زیرا فروشگاههای ۲ و ۳ حالتی شبیه بر حسب ویژگی های E, Q و L دارند اما فروشگاه ۲ سود داشته است در حالیکه فروشگاه ۳ ضرر داشته است. بنابراین اطلاعات در جدول برای پاسخ دادن به این سوال کافی نیست در نتیجه میتوان یک پاسخ ناقص به این سوال داد. می توان در نظر گرفت که اگر یک فروشگاه مشخص دارای ارزش high برای ویژگی E باشد آن فروشگاه دارای سود خواهد بود و اگر ارزش ویژگی E کم باشد (low) آن فروشگاه ضرر خواهد داد. بنابراین به طور مطمئن می توانیم بگوییم که فروشگاه ۱ و ۶ دارای سود هستند. فروشگاههای ۴ و ۵ ضرر می دهند و در مورد فروشگاههای ۲ و ۳ نمی توان گفت که سود دارند یا ضرر می دهند. با بکارگیری ویژگیهای E, Q و L می توان گفت که فروشگاههای ۱ و ۶ مطمئنا دارای سود هستند و ازه مطمئنا متعلق به مجموعه {1,3,6} است. در حالیکه فروشگاههای 1,2,3,6 ممکن است دارای سود باشند و ازه شاید (ممکن است) متعلق به مجموعه {1,3,6} باشند. در نتیجه می توان گفت که مجموعه {1,6} تقریب پایینی از مجموعه {1,3,6} و مجموعه {1,2,3,6} تقریب بالایی از مجموعه {1,3,6}

می باشد. مجموعه {2,3} تفاوت بین تقریب بالایی و پایینی است و به عنوان ناحیه مرزی از مجموعه {1,3,6} در نظر گرفته می شود.

در هر سیستم اطلاعاتی ویژگیها را می توان به دو دسته تقسیم کرد:

(۱) ویژگیهای تصمیم

(۲) ویژگیهای موقعیت

در جدول تصمیم مربوط به فروشگاهها ویژگی P به عنوان یک ویژگی تصمیم است در حالیکه ویژگیهای E و Q و L به عنوان ویژگیهای موقعیت شناخته می شوند. ویژگیهای تصمیم معین می کنند که چه تصمیم هایی باید اجرا شود تا ویژگیهای موقعیت بصورت متقاعد در آیند. در هر سطر از یک جدول تصمیم می توان قواعد تصمیم را به فرم اگر..... آنگاه..... بدست آورد. به عنوان مثال در جدول تصمیم ارائه شده داریم:

if (E=high)and (Q=good)and(L=no),then(P=profit)

if (E=no)and (Q=ave)and(L=no),then(P=loss)

قواعد متناقض (مخالف): قواعدی هستند که دارای ویژگیهای موقعیت یکسان اما ویژگیهای تصمیم متفاوت می باشند.

قواعد سازگار (موافق): قواعدی هستند که دارای ویژگیهای موقعیت یکسان و ویژگیهای تصمیم یکسان می باشند.

در مثال ارائه شده (فروشگاهها) قواعد ۲ و ۳ جزء قواعد متناقض و بقیه قواعد از قواعد سازگار می باشند. قواعد سازگار تحلیلگر را قادر می سازد تا یک تصمیم منحصر به فرد بگیرد ولی در مورد قواعد متناقض نمیتواند تصمیم یکتایی را در نظر بگیرد.

حد پایداری (اندازه استحکام) جدول تصمیم:

نسبت قواعد سازگار به همه قواعد را در یک جدول تصمیم حد پایداری گویند. اگر حد پایداری یک جدول تصمیم یک باشد (یعنی تمام قواعد جدول تصمیم سازگار باشند) در نتیجه می توان به صورت منحصر به فرد تصمیم گیری کرد ولی اگر حد پایداری یک سیستم اطلاعاتی کمتر از یک باشد می توان یک تصمیم نسبتاً کافی را در نظر گرفت. تئوری راف ست ابزاری را برای تولید، تحلیل و بهینه سازی مجموعه قواعد تصمیم از جدول داده ها را فراهم می کند. جدول تصمیم می تواند به عنوان یک نتیجه از مشاهدات (اندازه گیری)، شبیه سازی کامپیوتری و معرفت (دانش) افراد خبره بدست آید. چندین سیستم نرم افزاری مبتنی بر تئوری راف به شکل زیر می باشد:

LEERS Rough Das Rough Class DATALOGIC

برای مشخص کردن یک فرایند کنترلی جدول تصمیم می تواند مفید باشد. بنابراین قواعد تصمیم می تواند به عنوان قواعد کنترلی مورد استفاده قرار گیرد. در نتیجه تئوری مجموعه راف می تواند برای نتیجه گیری یک مجموعه بهینه از قواعد کنترلی از جدول سیستم به کار گرفته شود.

تئوری مجموعه راف در تحلیل قواعد اگر..... آنگاه..... [7]:

فرض کنید که U یک مجموعه جهانی از اشیاء باشد:

$$U = \{x, y, \dots\}$$

$$T = \{A_1, A_2, \dots, A_n\}$$

T مجموعه ویژگیهاست:

$$\text{Dom}(A_i)$$

مجموعه ای از ارزشهای ویژگی A_i :

$$\text{Dom} = \text{dom}(A_1) \cup \text{dom}(A_2) \cup \dots \cup \text{dom}(A_n)$$

هر عضو در مجموعه U را می توان به صورت منحصر به فرد به صورت یک نگاشت نمایش داد:

$$A \in T \text{ برای } t: T \rightarrow \text{Dom} \quad t(A) \in \text{dom}$$

هر جدول اطلاعاتی شامل موارد زیر می باشد:

- 1) $U = \{u, v, \dots\}$
- 2) $T = \{A_1, A_2, \dots, A_n\}$
- 3) $\text{Dom}(A_i)$
 $\text{Dom} = \text{dom}(A_1) \cup \text{dom}(A_2) \cup \dots \cup \text{dom}(A_n)$
- 4) $\rho: U \times T \rightarrow \text{Dom}$ تابع توصیف

تابع توصیف یک نگاشت است که:

$A_i \in T$ و $u \in U$ برای $\rho(u, A_i) \in \text{dom}(A_i)$

باید در نظر داشت که تابع توصیف ρ باعث یک مجموعه از نگاشتها می شود و هر نگاشت را یک چند تایی می نامند:
 $t = (\rho(u, A_1), \rho(u, A_2), \dots, \rho(u, A_i), \dots, \rho(u, A_n))$

توجه کنید که چند تایی در یک جدول اطلاعات الزاما وابسته به شیء (موضوع) منحصر به فردی نیست و دو شیء متفاوت می توانند دارای نمایش چند تایی یکسان باشند که این در مورد پایگاه داده های وابسته غیر مجاز می باشد. جدول تصمیم (DECISION TABLE) را با مشخصه (U, T, V, ρ) نمایش می دهند که در آن U مجموعه جهانی و T مجموعه ویژگیهاست که از اجتماع دو مجموعه غیر تهی C و D تشکیل شده است.
 $T = C \cup D$
 C مجموعه ویژگیهای موقعیت و D مجموعه ویژگیهای تصمیم می باشد.

وابستگی های دانش:

باتوجه به قوانین ریاضی یک کلاسه سازی یا پارتیشن بندی یک تجزیه از حوزه U به زیر مجموعه های گسسته می باشد که به آن کلاسه های هم ارزی نیز می گویند. هر پارتیشن توسط یک رابطه باینری تعریف می شود که این رابطه باینری یک رابطه هم ارزی است و سه خاصیت انعکاسی (بازتابی)، متقارن و انتقالی (تراگذر) را دارا می باشد

ویژگیها و روابط هم ارزی:

در جدول اطلاعات (داده ها) هر زیر مجموعه از ویژگیها یک رابطه هم ارزی فراهم می کند. اگر مجموعه B یک زیر مجموعه غیر تهی از T و دو شیء (موضوع) u و v (غیر قابل تشخیص توسط B) متعلق به U باشند خواهیم داشت:

$$u \equiv v \pmod{B} \text{ if } \rho(u, A_i) = \rho(v, A_i) \text{ برای هر ویژگی } A_i \text{ در } B \text{ داریم:}$$

علامت \equiv نشان دهنده رابطه هم ارزی است و به آن رابطه تعادل نیز می گویند و آن را با $\text{IND}(B)$ نمایش می دهند. کلاسه های هم ارزی شامل U را بوسیله $\text{IND}(B)$ یا ساده تر به صورت $[u]_B$ نشان می دهند که B می تواند یک مجموعه منفرد باشد.

$$B = \{A_i\}$$

در این حالت به سادگی رابطه هم ارزی به صورت $\text{IND}(A_i)$ را تعریف می کنیم و خواهیم داشت:

$$\text{IND}(B) = \bigcap \{ \text{IND}(A_i) : A_i \text{ in } B \}$$

رابطه بالا بیان می کند که کلاسه های هم ارزی مربوط به $\text{IND}(B)$ همه اشتراکهای ممکن از کلاسه های هم ارزی مربوط به $\text{IND}(A_i)$ را شامل می شود. با استفاده از مجموعه B که زیر مجموعه غیر تهی از T است می توان یک مجموعه از رابطه های هم ارزی که آن را با $\text{RCol}(B)$ نشان می دهند بدست آورد و آن را به صورت زیر تعریف می کنیم:

$$\text{Rcol}(B) = \{ \text{IND}(A_i) : A_i \in B \}$$

اگر مجموعه مورد نظر مجموعه ای از همه ویژگیها باشد در آن صورت مجموعه متشکل از رابطه های هم ارزی به صورت زیر خواهد بود:

$$\text{Rcol}(T) = \{ \text{IND}(A_1), \text{IND}(A_2), \dots, \text{IND}(A_3), \dots, \text{IND}(A_n) \}$$

پایه های دانش پاولاک (pawlak):

اگر Rcol یک مجموعه از رابطه های هم ارزی باشد، بر اساس تعریف زوج مرتب زیر را پایه دانش پاولاک نامیده اند.

$$K=(U,Rcol)$$

پروفسور زدیسلو و پاولاک مجموعه رابطه های هم ارزی را دانش (knowledge) نامید زیرا دانش ما درباره یک حوزه اغلب بوسیله کلاس بندیها نمایش داده می شود.

P و Q دو رابطه هم ارزی در K را در نظر بگیرید:

اگر کلاسهای هم ارزی Q یک اجتماع از کلاسهای هم ارزی P باشند در آن صورت می گوئیم Q وابسته به P می باشد یا اصطلاحاً Q سخت تر از P یا P نرمتر از Q است. به این ارتباط، وابستگی دانش (Knowledge Dependency) یا KD گویند.

بدیهی است که اشتراک دو رابطه هم ارزی یک رابطه هم ارزی دیگری است که شامل همه کلاسهای هم ارزی مشترک بین دو رابطه هم ارزی است به طور کلی اشتراک همه رابطه های هم ارزی بوسیله IND(Rcol) که به صورت یک رابطه هم ارزی دیگر است مشخص می شود. همچنین IND(Rcol) به عنوان یک رابطه علی-معلولی در RCol شناخته می شود.

اگر PCol و QCol دو زیر مجموعه از RCol باشند بنا به تعریف داریم:

$$1) QCol \text{ وابسته به } PCol \text{ است اگر } IND(Qcol) \text{ سخت تر از } IND(Pcol) \text{ باشد.}$$

2) PCol و QCol هم ارز یکدیگر هستند اگر:

$$\text{if } Pcol \text{ then } Qcol \quad , \quad \text{if } Qcol \text{ then } Pcol$$

یا به شکل ساده تر PCol و QCol هم ارز یکدیگر هستند اگر $IND(Pcol)=IND(Qcol)$

$$3) PCol \text{ و } QCol \text{ مستقل هستند اگر نه } \text{if } Pcol \text{ then } QCol \text{ و نه } \text{if } Qcol \text{ then } Pcol$$

در نتیجه ما می توانیم یک جدول اطلاعات را به عنوان یک پایه دانش پاولاک $(U, Rcol(T))$ در نظر بگیریم و مفهوم وابستگی دانش را برای جدول اطلاعات به کار ببریم.

قضیه: یک تناظر یک به یک جدول اطلاعات پایه های دانش پاولاک وجود دارد.

وابستگی تابعی:

هنگامیکه ارزشهای یک چندتایی روی مجموعه ای از ویژگیهای منحصر به فرد تعیین کننده ارزشهای چندتایی دیگر مجموعه از ویژگیها باشد گوئیم یک وابستگی تابعی اتفاق افتاده است.

جدولها و قواعد تصمیم:

تئوری مجموعه راف یک روش مؤثر برای استخراج قواعد از جدولهای اطلاعات و یا صریح تر از جداول تصمیم می باشد.

الف: candidate key:

در پایگاههای اطلاعاتی وابسته یک مجموعه از ویژگیها مانند K (مجموعه مینیمال) یک candidate key نامیده می شود اگر همه ویژگیها یک وابسته تابعی روی K باشند.

Extentional Candidate key یک روش خاص از تقلیل یا زدودن می باشد.

اگر $S=(U, T=CUD)$ یک جدول تصمیم باشد به طوریکه:

$$C=\{A1, A2, \dots, Ai, \dots, An\}$$

$$D=\{B1, B2, \dots, Bi, \dots, Bn\}$$

که C ویژگیهای موقعیت و D ویژگیهای تصمیم می باشد.

می توان گفت که دو پایه دانش روی مجموعه U به صورت زیر وجود دارد:

$$Rcol(C) = \{IND(A1), \dots, IND(A2), \dots, IND(An)\}$$

$$Rcol(D) = \{IND(B1), \dots, IND(B2), \dots, IND(Bn)\}$$

S یک جدول تصمیم سازگار است اگر داشته باشیم:

اگر $Rcol(C)$ آنگاه $Rcol(D)$

B رامی توان یک تقلیل از S در نظر گرفت اگر B یک زیر مجموعه مینیمال از C باشد و داشته باشیم:

اگر $Rcol(B)$ آنگاه $Rcol(D)$

بدیهی است که مجموعه B لازم نیست به صورت منحصر به فرد باشد.

اگر ما مجموعه D را مساوی مجموعه T انتخاب کنیم (D=T) در آن صورت یک تقلیل Extentional Candidate key خواهیم داشت.

ب: تئوری مجموعه راف یک روش مؤثر برای استخراج قواعد از جداول اطلاعاتی (IT) یا با صراحت بیشتر از جداول تصمیم (DT) می باشد. در هر جدول اطلاعاتی می توان چندین جدول تصمیم را مشاهده نمود. هر شیء (موضوع) در U در DT می تواند به عنوان یک قاعده تصمیم گیری تفسیر شود. متخصصین مجموعه راف علاقه مند به ساده کردن قواعد تصمیم هستند.

مثال: سیستم اطلاعاتی (جدول اطلاعاتی) زیر را در نظر بگیرید:

U	LOCATION	TEST	NEW	CASE	RESULT
ID-1	Houston	10	92	03	10
ID-2	San Jose	10	92	03	10
ID-3	Palto Alto	10	90	02	10
ID-4	Berkeley	11	91	04	50
ID-5	Newyork	11	91	04	50
ID-6	Atlanta	20	93	70	99
ID-7	Chicago	20	93	70	99
ID-8	Baltimore	20	93	70	99
ID-9	Seattle	20	93	70	99
ID-10	Chicago	51	95	70	94
ID-11	Chicago	51	95	70	95

A: یک جدول تصمیم (DT) انتخاب می کنیم و آن را S می نامیم در نتیجه خواهیم داشت: $S=(U, T=CUD)$

$C=\{TEST, NEW, CASE\}$ $D=\{RESULT\}$

U: مجموعه جهانی D: مجموعه ویژگیهای تصمیم

C: مجموعه ویژگیهای موقعیت

B: فرمهای ID-10 و ID-11 دو قاعده ناسازگارند و بقیه فرمها قاعده سازگار هستند. در جداول تصمیم قاعده هایی را در نظریه گیریم که سازگار باشند.

C: با استفاده از رابطه هم ارزی IND(RESET) می توان اشیاء (ID-1 تا ID-9) را در سه کلاس هم ارزی که کلاسهای تصمیم نامیده می شوند دسته بندی نمود.

DECISION1={ID-1,ID-2,ID-3}=[10] RESULT,
DECISION2={ID-4,ID-5}=[50] RESULT,
DECISION3={ID-6,ID-7,ID-8,ID-9}=[99] RESULT

D : با استفاده از رابطه هم ارزی ©IND می توان اشیاء را در چهار کلاس هم ارزی که کلاسهای موقعیت نامیده می شوند تقسیم بندی نمود.

CASE1={ID-1,ID-2} CASE2={ID-3} CASE3={ID-4,ID-5}
CASE4={ID-6,ID-7,ID-8,ID-9}

E: وابستگی دانش (KD):

بعضی از اشیاء توسط ویژگیهای موقعیت و تصمیم غیر قابل تشخیص می باشند.

CASE1 زیر مجموعه ای از DECISION1 است.

CASE2 زیر مجموعه ای از DECISION1 است.

CASE3 زیر مجموعه ای از DECISION2 است.

CASE4 زیر مجموعه ای از DECISION3 است.

عبارتهای فوق دلالت دارد بر اینکه رابطه هم ارزی IND(RESET) وابسته به IND(C) است یا به طور معادل RESULT یک KD روی C می باشد.

F : استنتاج قواعد:

با استفاده از روابط بالامی توان قواعد زیر را استنتاج کرد:

IF TEST=10,NEW=92,CASE=03 THEN RESULT=10
IF TEST=10,NEW=90,CASE=02 THEN RESULT=10
IF TEST=11,NEW=91,CASE=04 THEN RESULT=50
IF TEST=20,NEW=93,CASE=70 THEN RESULT=99

G: تقلیل ها:

واضح است که {TEST,NEW} و {CASE} دو تقلیل هستند. بنابراین موقعیتها روی ویژگی CASE می توانند از قواعد داده شده در بالا حذف شوند.

H: تقلیل کننده های ارزشی:

توجه کنید که قاعده ۱ از رابطه (CASE1 زیر مجموعه DECISION1 است) مشتق شده است پس ما می توانیم توصیف از CASE1 را به وسیله قید TEST=10 به تنهایی نشان دهیم که آن تقلیل کننده ارزشی نامیده می شود.

CASE1={U:u.TEST=10}

با استفاده از استدلال مشابه می توان چهار قاعده را به صورت زیر خلاصه نمود.

1. IF TEST=10 THEN RESULT=10
2. IF TEST=11 THEN RESULT=50
3. IF TEST=20 THEN RESULT=99

این تقلیل منحصر به فرد نمی باشد و می توان با انتخابهای متفاوت تقلیلهای متفاوت به دست آورد.

مزایای کاربرد تئوری راف:

تئوری راف کاربردهای زیادی در مهندسی، تحلیل داده های پزشکی، پردازش تصویر و... دارد. برخی از مزایای کاربردی تئوری مجموعه های راف به صورت زیر می باشد [4]:

* یک الگوریتم موثر برای یافتن الگوهای پنهان در داده ها

* یافتن مجموعه های مینیمال داده ها (کاهش یا تقلیل داده ها)

* ارزیابی اهمیت داده ها

* تولید مجموعه های مینیمال از قواعد تصمیم گیری از داده ها

* سادگی فهم و تفسیر آسان نتایج الگوریتم

مثالهای کاربردی تئوری مجموعه راف:

* استخراج قواعد با استفاده از مجموعه راف هنگامیکه ارزشهای تابع تعلق به صورت فواصلی می باشند.

* کاربرد مجموعه راف با کشف کننده ها برای انتخاب خصوصیت (طرح)

* روشی برای استخراج تغذیه های پایگاه داده ها با استفاده از مجموعه راف

* کاربرد تئوری راف در کنترل (کنترل راف) [6]

* کاربرد مجموعه راف برای داده کاوی در سیستم های اطلاعاتی بیمارستان

* تحلیل عملی روی مجموعه اطلاعات مراقبت ژنتیک با استفاده از مجموعه راف

* تولید اتوماتیک داستان با روشهای NLP و مجموعه های راف [8]

* کلاسه سازی نقطه بهره برداری سیستم قدرت [9]

* تشخیص خطا در سیستم های قدرت [10]

* استخراج معرفت از پست های توزیع [11]

* استخراج قواعد با تحلیل داده های پزشکی [12]

* کاهش (فشرده سازی) معرفت در جداول اطلاعاتی (اطلاعات خودروها) [13]

* تشخیص خطا در سیستم های حمل و نقل [14]

* کار بر روی هوش پیوندی با استفاده از مجموعه راف

* استخراج معرفت جاسازی شده در شبکه عصبی تعلیم یافته با استفاده از مجموعه راف

* کشف روابط بین ویژگیها و وابستگی قواعد با استفاده از مجموعه راف

* استدلال درباره معرفت با استفاده از مجموعه راف

* مدل سازی تحلیلی سیستم دما با استفاده از مجموعه راف

* مدل سازی مجموعه بیماران قلبی موروئی با استفاده از مجموعه راف

* توسعه چک لیست تشخیصی در یک اتاق اورژانس با استفاده از مجموعه راف

* کلاسه سازی کلمات دور افتاده (پرت) زبان رسمی کشور تایلند با استفاده از مجموعه راف

* بدست آوردن قواعد از اطلاعات ناقص و متناقض با استفاده از مجموعه راف

* روشی برای ارزیابی کیفیت در سیستم حمل و نقل با استفاده از مجموعه راف

* استفاده از مجموعه راف برای انطباق رهنمودهای بالینی برای سازمانهای مراقب سلامتی

نتیجه:

در این مقاله، مفاهیم اولیه و زمینه های تئوری و کاربردی مجموعه های راف در استخراج دانش (معرفت) و داده کاوی و کاربردهای عملی و واقعی آن در زمینه های مختلف مهندسی و علوم دیگر معرفی شد. همچنین به صورت مختصر تئوری مجموعه های راف در تحلیل قواعد اگر.....آنگاه.....همراه با ذکر مثال بیان شد که می توان مباحث تئوری بیشتر را درباره مجموعه های راف در منابع موجود یافت. با توجه به کارهای انجام شده درباره مجموعه های راف می توان نتیجه گرفت که این تئوری یک ابزار قدرتمند در بازیابی اطلاعات، داده کاوی [15]، تقلیل (کاستن) داده های زائد از پایگاه داده ها می باشد.

مراجع:

- [1] Honby, A.S., "Oxford Advanced Learners Dictionary of Current English", Oxford University Press, UK, 1974
- [2] Ziarko, W., "The Discovery ,Analysis and Representation of Data Dependencies in Databases", Knowledge Discovery in Databases, AAAI MIT Press ,Cambridge ,MA, 1993, pp.213-228
- [3] Pawlak, Z., "Rough Sets", International Journal of Computer and Information Sciences, Vol.11, 1982, pp.341-356
- [4] Pawlak, Z., "Rough Sets and Data Analysis" Proceeding of IEEE Conference ,ISSN:0-7803-3687-9, 1996
- [5] Pawlak, Z., "Rough Sets : Theoretical Aspects of Reasoning about Data ", Kluwer Academic Publishers, Dordrecht ,boston ,London, 1991.
- [6] Toshinori, M., "Rough Control Application of Rough Set Theory to Control", Computer and Information Science Department Cleveland State University
- [7] T.Y.Lin, "An Overview of Rough Set Theory from the Point of View of Relational Databases", Department of Electrical Engineering and Computer Science, University of California
- [8] Beaubouef, T., Lang, R., "Rough Set Techniques for Uncertainty Management in Automated Story Generation" Proceeding of the 36th Annual Conference on Southeast Regional Conference, April, 1998, pp.326-331
- [9] Fereidunian, A.R., Lesani, H., Lucas, C., "Distribution Systems Reconfiguration Using Pattern Recognizer Neural Networks", International Journal of Engineering (IJE), Transactions B: Applications, Vol. 15, No.2, pp.135-144, July 2002
- [10] Zhang ,Q., Han, Z., Wen, F., "A New Approach for Fault Diagnosis in Power Systems Based on Rough Set Theory ", Proceeding of the 4th international Conference on Power Systems Control ,Operation and Management, APSCOM-97, Hong Kong, November 1997, pp.597-602
- [11] Hor, C.-L., Crossley, P., Dunand, F., "Knowledge Extraction within Distribution Substation Using Rough Set Approach ", IEEE Power Engineering Society Winter Meeting 2002 Proceeding 0-7803-7322-7/02, January 27-31 ,2002 , New York, NY, USA
- [12] Nakayama ,H., Hattori, Y., Ishii, R., "Rule Extraction Based on Rough Set Theory and Its Application to Medical Data Analysis", IEEE Conference Proceeding 0-7803-5731-0/99, 1999, pp.v924-v929
- [13] Douqian, M., Jue, W., "Information-Based Algorithm for Reduction of Knowledge ", Proceeding of IEEE International Conference on Intelligent Processing Systems, Beijing, China, October 28-31, 1997.
- [14] Xiaolei, L., Xiaobing, W., "The Application of Rough Set Theory in Vehicle Transmission System Fault Diagnosis ", Proceeding ISSN:0-7803-5296-3/99, 1999, pp.240-242
- [15] Grzymala-Busse, J., Ziarko, W., "Data Mining and Rough Set Theory", Communications of ACM, April 2000, Vol.43, No.4, pp.108.