



ارائه شده توسط:

سایت ترجمه فا

مرجع جدیدترین مقالات ترجمه شده

از نشریات معتبر

یک طبقه بندی بیز ساده برای مدارک علمی

چکیده

طبقه‌بندی دستی افراد به دسته‌های مختلف براساس مدرک تحصیلی‌شان، کار خسته‌کننده‌ای است و این ممکن است که تغییری نسبت به سناریوی بحث شده، ایجاد کند. این مقاله یک روش طبقه‌بندی با استفاده از معیار الگوریتم طبقه‌بندی بیز ساده برای طبقه‌بندی افراد به دسته‌های مختلف براساس نمایش برخی ویژگی‌های مدارک علمی‌شان پیشنهاد می‌کند. نتایج آزمایشی قابل ارزیابی نشان می‌دهد که روش طبقه‌بندی پیشنهادی می‌تواند امیدوارکننده باشد و در جای دیگر استفاده شود. روش پیشنهادی، به صورت آزمایشی با 90٪ دقت مقادیر κ ، اثبات کارایی بالای آن را تأیید می‌کند. این روش طبقه‌بندی می‌تواند کار دستی این جهان را کاهش دهد و به راحتی در دسته‌بندی کمک کند.

کلمات کلیدی: طبقه‌بندی، داده‌کاوی، مدارک علمی، κ ، بیز ساده

1. مقدمه

متعدد پیش آمده که یک فرد براساس تجزیه و تحلیل مدارک علمی که در زندگی کسب کرده است، در نظر گرفته شود. در چنین مواردی، طبقه‌بندی افراد با توجه به مدرک تحصیلی آموزشی‌شان می‌تواند به تصمیمی فنی آزادانه، و عاری از هرگونه تعصب بسیار کمک کند و از این رو می‌تواند قابل اجرا باشد.

این مقاله روشی برای دسته‌بندی مدارک علمی با استفاده از معیار الگوریتم طبقه‌بندی بیز ساده پیشنهاد می‌کند. این روش می‌تواند در بسیاری از برنامه‌های کاربردی مانند تفکیک براساس فهرست کوتاه ارتباط آموزشی، برای استخدام افراد براساس میزان تحصیلات و غیره، استفاده شود. سازمان این مقاله در زیر آورده شده است: بخش 2 شامل بررسی ادبیات می‌شود. بخش 3 الگوریتم بیزین ساده و روش طبقه‌بندی پیشنهادی را توضیح می‌دهد. بخش 4 تجزیه و تحلیل نتایج آزمایشی براساس جدول بندی ذکر شده است و بخش 5 مقاله را نتیجه‌گیری می‌کند.

الگوریتم بیزین ساده یک الگوریتم طبقه‌بندی کلاسیک است که بهره‌وری خود را در برنامه‌های مختلف و چند مقاله نمایش بهره‌وری از طبقه‌بندی کننده که در اینجا ورد بحث است، ثابت کرده است.

مقاله Mauricio A. Valle و همکاران روش پیش‌بینی ویژگی‌های تعیین‌کننده در مورد یک الگوریتم طبقه‌بندی بیز ساده شامل یک روش تست براساس اعتبارسنجی متقابل را بحث می‌کند. تأیید آزمایشی صفات اجتماعی و جمعیت شناختی است که به پیش‌بینی عملکرد آینده عامل فروش در یک مرکز تلفن کمک نمی‌کند.

Dunja Mladenic و همکاران. پژوهش با انتخاب ویژگی‌های کمک برای طبقه‌بندی با استفاده از مشخصات خاص و توانایی یادگیری طبقه‌بندی بیش از یک داده متنی که توزیع نابرابری است. زمانی که دامنه و ویژگی‌های الگوریتم طبقه‌بندی شده، در نظر گرفته شود، عملکرد طبقه‌بندی افزایش می‌یابد. Dong Tao و همکاران. مقاله بهبود الگوریتم بیزی ساده را با ترکیب روش کلاسیک با یک روش انتخاب ویژگی براساس شاخص Gini پیشنهاد می‌کند. این روش ترکیبی، عملکرد طبقه‌بندی متن را بهبود می‌بخشد.

Kabir Md Faisal و همکاران. پژوهش ترکیب روش خوشه‌بندی k-means با الگوریتم طبقه‌بندی ساده برای افزایش دقت. این روش خوشه‌بندی گروه‌های نمونه‌های آموزشی را به دسته‌بندی‌های مشابه، پس از آن همه گروه‌ها تحت طبقه‌بندی بیز ساده آموزش داده می‌شوند. این روش برای تأیید بهبود دقت است.

Santra A.K. و همکاران. تحقیقات ثابت می‌کند که در مورد استفاده از کاربرد وب، در حالی که از یک طبقه‌بندی بیز ساده به جای درخت تصمیم استفاده می‌کند، که زمان صرف شده برای طبقه‌بندی و حافظه کاهش پیدا می‌کند. مقاله نشان می‌دهد که ماهیت استقلال شرطی ویژگی‌ها در الگوریتم اصلی بیز ساده در بعضی موارد ضعیف به نظر می‌رسد و روش وزندهی محلی که از الگوریتم کلاسیک از نظر دقت بهتر است، پیشنهاد می‌کند. Pradeepta K. Sarangiet و همکاران. مقاله استخراج ویژگی با استفاده از تجزیه LU به دنبال استفاده از طبقه‌بندی بیز ساده برای تشخیص الگو توصیف می‌کند. این کاربرد جهانی طبقه‌بندی را نشان می‌دهد.

Yildirim و Birant D. مقاله تحقیقاتی تأیید آزمایشی اثر توزیع‌های مختلف بر روی ویژگی‌های مورد بحث. مشاهده شده است که کاربرد توزیع براساس طبیعت ویژگی‌ها به جای استفاده از یک توزیع در سراسر تمام ویژگی‌ها دقت را افزایش می‌دهد. AberBadr El Din Ahmed و Ibrahim Sayed Elarabarticle بحث در مورد استفاده از الگوریتم‌های طبقه‌بندی برای پیش‌بینی نمره نهایی دانش آموزان است. مقاله Ron Kohavi ترکیب طبقه‌بندی بیز را با درخت تصمیم که به عنوان NBTree نامیده است، برای افزایش دقت طبقه‌بندی پیشنهاد می‌کند.

همچنین دریافتند که استقلال شرطی کلاس در مورد مجموعه داده کوچک منفعل است اما در صورت مجموعه داده‌های بزرگ، این فرض منجر به اشکال در طبقه‌بندی و کاهش دقت می‌شود.

Shasha Wang و همکاران. مقاله نسخه ارتقا یافته طبقه‌بندی NBTree ترکیبی را پیشنهاد و آن را به‌عنوان NBTree چندگانه به نام MNBTree، که در آن یک طبقه‌بند بیزین ساده چندجمله‌ای برای گره‌های برگ درخت تصمیم کاربرد دارد. علاوه‌براین، برای افزایش عملکرد، بدهاۀ دیگری با گنجاندن طبقه‌بندی چند کلاسه ساخته شده است و سیستم به‌عنوان نسخه چند کلاسه MNBTree نامیده می‌شود.

با توجه به مقالات پژوهشی فوق‌الذکر، جوانب مثبت الگوریتم طبقه‌بندی بیزین به‌طور کامل مطالعه شده و دریافتند که این الگوریتم بهترین خواهد شد با توجه به ماهیت داده‌های مورد استفاده برای آزمایش متشکل از هردو داده عددی و متنی که به‌طور مستقل کمک به طبقه‌بندی باشد.

2. سیستم پیشنهادی

2.1. مرور کلی بیز ساده

الگوریتم طبقه‌بندی بیز بسیار ساده است که در آن فرض بر این است که ویژگی‌های طبقه‌بندی مستقل هستند و بین آنها هیچ همبستگی وجود ندارد. بسیاری از محققان دریافته‌اند که این فرض استقلال در تمام موارد که دیگر روش‌های جایگزین ارائه شده، برای افزایش عملکرد کار نمی‌کند. یکی از روش‌های جایگزین را Liang xiao Jiang ارائه داده است.

تکنیک اصلی بیز ساده بر احتمال شرطی و حداکثر وقوع احتمال است. الگوریتم بیز ساده براساس توضیحات ارائه شده به شرح زیر است:

- فرض کنید G مجموعه آموزش با N متغیر است که در آن K به‌عنوان بردار ویژگی بُعد M در $M = \{M_1, M_2, \dots, M_k\}$ نشان می‌دهد.

- برای 'P' کلاس‌های C_1, C_2, \dots, C_p وجود دارد. با توجه به این طبقه‌بندی بیز ساده، T یک متغیر متعلق به کلاس C_x است که فقط زمانی آن را به یک احتمال شرطی بالاتر از هر کلاس دیگر C_y می‌برد، که در آن $x \neq y$ باشد.

$$P(C_x | T) > P(C_y | T) \text{ و } P(C_x | T) = (P(T | C_x) * P(C_x)) / P(T)$$

- از کلاس مستقل مشروط فرض شده است،

$$P(M_k | C_x) = P(M_1) * P(M_2) * P(M_3) \dots * P(M_k) \quad P(M | C_x) = \prod_{i=1}^k P(M_i | C_x)$$

- کلاس C_x به عنوان طبقه خروجی پیش‌بینی شده $P(M | C_x) * P(C_x) > P(M | C_y) * P(C_y)$ که در آن x, y و $1 \leq x, y$ است.

2.2 شرح دیتاست

سناریویی که اینجا مطرح شده با جزئیات مدارک علمی افراد که آنها در دوره آموزشی‌شان کسب کرده‌اند سروکار دارد. براساس این جزئیات، افراد به کلاس‌های مختلف طبقه‌بندی می‌شوند. در اینجا سه کلاس در نظر گرفته شده و آنها عبارتند از:

- کم.

- متوسط.

- زیاد.

ویژگی‌های گرفته شده در دست بررسی در زیر آورده شده:

- عدد کل سال‌هایی که برای تحصیل صرف شده.

- درجه تحصیلی کسب شده.

- موضوع اصلی / رشته تخصصی مربوط به سناریو تحت بررسی.

عدد کل سال‌هایی که برای آموزش سپری شده به عنوان یک ویژگی تعیین عمق دانش به دست آمده، به طور مستقیم به زمان سرمایه‌گذاری در یادگیری متناسب با آموزش در نظر گرفته شده است.

درجه علمی به دست آمده توسط یک فرد، نقش حیاتی را در تصمیم‌گیری مهارت شخص بازی می‌کند.

زمینه‌های تخصصی افراد به عنوان ویژگی‌های در نظر گرفته در تعیین دلیل انعطاف‌پذیری این سیستم شده که در انواع حالات استفاده می‌شود و این ویژگی کاربرد وابسته است.

شامل مجموعه داده در حدود 25000 ورودی که در آن 20000 ورودی برای آموزش و باقی مانده 5000 ورودی

برای تست بهره‌وری طبقه‌بندی و تعیین چگونه خوب یاد گرفته، به طبقه‌بندی داده‌ها استفاده می‌شود. نمایش

مجموعه داده به صورت زیر است:

- نمایش کل داده‌ها عبارتست از $D = \{D_1, D_2, \dots, D_{25000}\}$
 - نمایش کل داده‌های آموزشی عبارتست از $\{D_1, D_2, \dots, D_{20000}\}$
 - نمایش داده‌های تست عبارتست از $\{D_{20001}, D_{20002}, \dots, D_{25000}\}$
- یک فرد با مدرک علمی مشابه ممکن است تحت تفکیک حالات مختلف شده باشد، چون طبقه‌بندی کاربرد وابسته است.

2.3 آموزش و تست مجموعه داده نمونه

این مجموعه آموزشی نمونه مطابق با انتخاب یک کاندید برای آموزش دپارتمان شیمی در زیر ارایه شده است. این جدول به‌وضوح ویژگی‌های مختلف مورد استفاده برای طبقه‌بندی و کاربرد وابسته را نشان می‌دهد. داده‌های تست مشابه داده‌های آموزش بدون ستون کلاس خواهد بود که می‌تواند با کمک اجرای الگوریتم زیر پیش‌بینی کند که در زیر ارائه شده است.

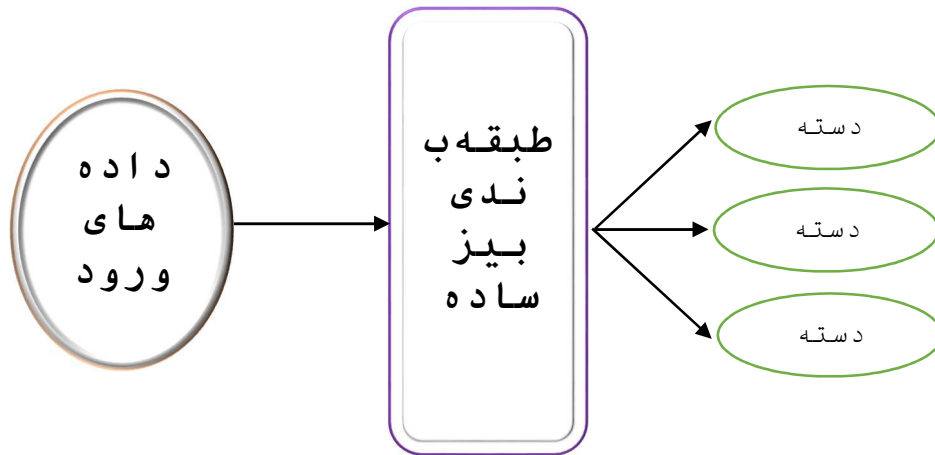
جدول 2.

مدرک تحصیلی	تعداد سال‌ها	تخصص	کلاس
کارشناسی علوم پایه	15	شیمی	متوسط
کارشناسی ارشد بازرگانی	17	تجارت	کم
کارشناسی ارشد علوم پایه	17	شیمی	زیاد
دکتر	25	شیمی	زیاد
کارشناسی ارشد مدیریت	18	مدیریت	کم
تجاری			
کارشناسی علوم پایه	15	فیزیک	کم
دکتر	22	زیست‌شناسی	کم

توجه: اگر یک فرد یا مدرک علمی بالا در درجه‌بندی PG، برچسب کلاس کم دارد به‌دلیل ویژگی مساحت تخصص کاربرد وابسته است.

2.4 روش طبقه‌بندی

داده‌های آموزشی برای دسته‌بندی داده‌ها به سه دسته تحت بررسی که در بخش بالا براساس ویژگی‌ها تعیین شده‌اند، استفاده می‌شوند. معماری برای چنین روش طبقه‌بندی در شکل 1 آورده شده است.



شکل 1. معماری طبقه‌بندی ارائه شده است.

2.5 الگوریتم

روش: الگوریتم بیز ساده برای مدارک علمی

- شروع

- مقداردهی اولیه

nc - تعداد دسته‌ها

na - تعداد ویژگی‌ها

N - تعداد نمونه‌ها

- برای هر کلاس C_i انجام می‌دهد

محاسبه احتمال قبلی $P(C_i) = \sum C_i / \sum N, i \in \{1, nc\}$

- برای هر کلاس C_i انجام می‌دهد

برای هر کلاس A_j انجام می‌دهد

محاسبه احتمال شرطی $P(A_j|C_i) = \sum C_i$

با $A_j / \sum C_i, i \in \{1, \dots, n_c\}$

و $j \in \{1, \dots, n_a\}$

• برای هر کلاس C_i انجام می‌دهد

محاسبه احتمال شرطی چند تایی K i.e $P(K|C_i) = P(A_1|C_i) * P(A_2|C_i) * \dots * P(A_{n_a}|C_i)$

• برای هر کلاس C_i انجام می‌دهد

محاسبه احتمال قبلی چند تایی K i.e $P(C_i) * P(K|C_i)$

• پیش‌بینی

اگر $((P(C_p) * P(K|C_p)) > (P(C_q) * P(K|C_q)))$

← پیش‌بینی C_p

در غیر این صورت

← پیش‌بینی C_q

که در آن $p, q \in \{1, \dots, n_c\}$ و $p \neq q$ باشد

• پایان

3. نتایج آزمایش و تجزیه و تحلیل

خروجی طبقه‌بندی تعداد نمونه‌های تست گرفته شده نشان می‌دهد که طبقه‌بندی در سه کلاس مختلف جدول است در جدول 1. براساس این جدول معیارهای عملکرد محاسبه شده است و در جدول 2 و جدول 3 فهرست شده‌اند. جدول 2 دقت و مقادیر طبقه‌بندی $kappa$ نشان می‌دهد. جدول 3 فهرست حساسیت، ویژگی و درجه توزیع سه سطح مختلف بحث شده را فهرست کرده است.

جدول 1. تعداد نمونه‌هایی که در سطح‌های مختلف دسته‌بندی شده‌اند.

پیش‌بینی /

دسته‌های

زیاد	متوسط	کم	انتظار
0	126	1233	کم
126	1716	100	متوسط
1566	133	0	زیاد

جدول 2. دقت و معیارهای Kappa طبقه‌بندی

مقدار	پارامترها
0.903	دقت
0.8528	Kappa

جدول 3. میزان حساسیت و گزینش‌پذیری برای سه سطح

دسته زیاد	دسته متوسط	دسته کم	پارامترها
0.9255	0.8689	0.9250	حساسیت
0.9598	0.9253	0.9656	ویژگی
0.3384	0.3950	0.2666	توزیع

3.1 معیارهای عملکرد

به‌طور سنتی، دقت، اصلی‌ترین معیار برای تعیین طبقه‌بندی‌کننده بود. اما وقتی که مجموعه داده به‌خاطر تمایزش به‌سمت کلاس غالب، کاملاً انحراف‌دار است، دقت، غیرقابل اعتماد است بنابراین مقادیر گوناگون دیگر نشان‌می‌دهند که می‌توان به‌وضوح توانایی حقیقی یک طبقه‌بندی‌کننده را توصیف کرد.

میزان کارایی برای تعیین عملکرد دسته‌بندی‌کننده‌هایی که در زیر، توضیح داده‌شده، استفاده می‌شود.

3.1.1 حساسیت:

حساسیت به معنای توانایی طبقه‌بندی‌کننده در شناسایی دسته‌ی مثبت از دسته‌ی مثبت و منفی و به‌طور صحیح از دسته‌ی منفی است. حساسیت، نسبت حقیقی (صحیح) ارزش‌های مثبت را نشان می‌دهد.

(منفی پیش‌بینی شده + مثبت پیش‌بینی شده) / مثبت پیش‌بینی شده = حساسیت

3.1.2 ویژگی:

ویژگی، توانایی طبقه‌بند را به‌طور مناسب در مستثنی نمودن ارزش دسته‌های مثبت از دسته‌های منفی و دسته‌های منفی از مثبت‌ها نشان می‌دهد. ویژگی، نسبت ارزش‌های حقیقی منفی را تشریح می‌کند.

(مثبت به اشتباه پیش‌بینی شده + منفی پیش‌بینی شده) / منفی پیش‌بینی شده = ویژگی

3.1.3 دقت:

دقت، ملاکی است برای تعیین مقداری از ارزش‌های دسته‌ی مثبت که به‌درستی از مثبت‌ها دسته‌بندی شده‌اند و مقداری از دسته‌ی منفی که دقیقاً به‌عنوان منفی علامت‌گذاری شده‌اند. دقت، نشان‌دهنده ارزش‌هایی است که به‌درستی دسته‌بندی شده‌اند.

(کل منفی‌ها + کل مثبت‌ها) / (منفی پیش‌بینی شده + مثبت پیش‌بینی شده) = دقت

3.1.4 Kappa:

ارزش Kappa نشان‌دهنده توافقی ارزیابی بینابین و شرایط مناسب توسط مقدار Kappa 1 است. به‌همان اندازه که ارزش کاهش می‌یابد، شرایط سطح نیز کاهش می‌یابد. Kappa، ارزش آماری نرمالی، که دقت سیستم تعریف‌شده توسط کاربر، و آنکه از یک سیستم فرضی تصادفی مقایسه می‌کند را نشان می‌دهد. این ثابت شده که بیشترین ملاک مؤثر در تعیین بازده طبقه‌بندی‌کننده، مقایسه سیستم‌شان با یک سیستم ایده‌آل است. این، همچنین تعیین می‌کند که طبقه‌بندی‌کننده چگونه دسته‌بندی را آموخته یا مقادیر آموزشی را به‌خاطر سپرده، به‌عبارت دیگر، طبقه‌بندی‌کننده رابطه بین ویژگی‌ها با روش دقیق را درک می‌کند و نه اینکه فقط کلاس مقادیر تکراری را، تکرار کند.

$$Kappa = (1 - \text{تصادفی}) / (\text{تصادفی} - \text{دقت})$$

تصادفی = ((منفی صحیح + مثبت ناصحیح) * (منفی صحیح + منفی ناصحیح) + (منفی ناصحیح + مثبت صحیح) * (منفی صحیح + مثبت صحیح))

* (مثبت ناصحیح + مثبت صحیح) * تقسیم‌بر (کل مثبت‌ها + کل منفی‌ها) * (کل مثبت‌ها + کل منفی‌ها)

3.1.5 توزیع

توزیع نسبت دسته مثبت به کل جمعیت آن است و به صورت

$$\text{توزیع} = \frac{\text{کل منفی ها} + \text{کل مثبت ها}}{\text{منفی به اشتباه پیش بینی شده} + \text{مثبت پیش بینی شده}}$$

بدست می آید.

3.1.6 نمادهای استفاده شده

جدول 4 شامل نمادهای مقادیر مورد استفاده برای محاسبه عملکرد ذکر شده بالا است. این مقادیر مبنای تعیین محاسبه عملکرد ذکر شده بالا است. این مقادیر معمولاً از ماتریس درهم ریختگی بدست آمده اند. یک ماتریس درهم ریختگی، یک ماتریس دو در دو دارای ارزش های اصلی کلاس مثبت و منفی است که توسط طبقه بند، بعد از آموزش، پیش بینی شده است.

جدول 4. نمادگذاری هایی که در محاسبه ملاک های عملکرد استفاده شده است:

مقداری که به درستی تشخیص داده شده (مثبت حقیقی)
مقداری که به اشتباه صحیح تشخیص داده شده (مثبت کاذب)
مقداری که به درستی ناصحیح تشخیص داده شده (منفی حقیقی)
مقداری که به اشتباه ناصحیح تشخیص داده شده (منفی کاذب)
مقادیر مشمول X به عنوان کلاس X شناخته شده و دسته Y به عنوان کلاس Y
مقادیر مستثنی Y به عنوان کلاس X شناخته نشده و دسته X به عنوان کلاس Y شناخته نشده
کل مثبت ها همه مقادیری که به عنوان دسته مثبت شناخته شده است
کل منفی ها همه مقادیری که به عنوان دسته منفی شناخته شده است

3.2 تحلیل:

کارایی نتایجی که الگوریتم طبقه بندی بیز ساده ارائه می کند، وقتی که مشخصات (نشانه ها)، غیر عددی هستند به خوبی اثبات می شود. چرا که اصالت طبیعی دسته بندی کننده، محاسبه براساس بیشترین درست نمایی حادثه

هستند، نتایج دسته‌بندی رضایت‌بخش‌اند و برخی دسته‌بندی‌کننده‌ها می‌توانند در کاربردهای گوناگون جایی که یک شخص می‌تواند براساس درجه‌ی علمی اکتسابی‌اش طبقه‌بندی کند، استفاده شود. دقت بالا و ارزش بالای Kappa نمایان می‌شود(اشاره دارد به) که دسته‌بندی توسط دانشی که درطول فرآیند آموزش، کسب‌شده، ساخته شده و از این‌رو ثابت می‌شود که یک سیستم امیدوارکننده است. ارزش بالای Kappa همچنین معرفی می‌کند آنچه را که سیستم به‌منظور طبقه‌بندی نسبت به حفظ مقادیر آموزش دیده است. ارزش بالای حساسیت و اختصاصی‌بودن، مشخص می‌کند که دسته‌بندی‌کننده، قادر به شناسایی یک کلاس(سطح)مثبت به‌صورتی که از مثبت‌ها و مانع‌شدن(مستثنی کردن) کلاس‌های منفی از مثبت و برعکس. الگوریتم ساده‌ی بیز که کارایی‌اش در وضعیت‌های گوناگون زندگی ثابت شده‌است، باردیگر کارایی‌اش در این سناریوی متعهدشده، اثبات شده‌است.

4. نتیجه‌گیری:

یک روش طبقه‌بندی براساس طبقه‌بندی بیز ساده برای طبقه‌بندی افراد مختلف در سطوح مختلف براساس ویژگی‌های مختلف مطابق با وضعیت تحصیلی‌شان، در این مقاله پیشنهاد می‌شود. کار آینده ممکن است با افزایش تعداد ویژگی‌های مربوط به تحصیل برای تعیین صلاحیت دانش اشخاص و دسته‌بندی آن‌ها برطبق آن، انجام‌شود. علاوه بر این، ممکن است الگوریتم‌های مختلف طبقه‌بندی دیگر، برای بررسی تغییرات در معیارهای عملکرد و تحلیل توانایی دسته‌بندی طبقه‌بندهای مختلف، مورد استفاده قرارگیرند.



این مقاله، از سری مقالات ترجمه شده رایگان سایت ترجمه فا میباشد که با فرمت PDF در اختیار شما عزیزان قرار گرفته است. در صورت تمایل میتوانید با کلیک بر روی دکمه های زیر از سایر مقالات نیز استفاده نمایید:

لیست مقالات ترجمه شده ✓

لیست مقالات ترجمه شده رایگان ✓

لیست جدیدترین مقالات انگلیسی ISI ✓

سایت ترجمه فا ؛ مرجع جدیدترین مقالات ترجمه شده از نشریات معتبر خارجی