

ارایه و بررسی روشی موثر جهت انتخاب صفات خاصه مناسب در ساخت درخت تصمیم

مهدی اسماعیلی^۱؛ منصور طرفدار^۲

چکیده

حذف و نادیده گرفتن صفات خاصه در پایگاه داده های حجیم جهت تصمیم گیری و تحلیل سریع می تواند موثر واقع شود. یکی از تکنیک های مهم و پر کاربرد در مراحل پیش پردازش داده ها، کاهش ابعاد داده ها و یا به عبارت دیگر، انتخاب صفات خاصه مناسب برای داده کاوی است.

در این مقاله تاثیر کاهش ابعاد در دسته بندی داده ها بررسی و پس از آن الگوریتمی موثر پیشنهاد می شود. الگوریتم از دو مرحله تشکیل شده است که در مرحله اول با کمک دو معیار، ترتیبی بین صفات خاصه به دست می آوریم. صفات خاصه مرتب شده در مرحله اول به عنوان ورودی برای مرحله بعدی که ساخت درخت تصمیم است استفاده می شود. نتایج نشان می دهند که درخت تصمیم حاصل از این الگوریتم پیشنهادی همراه با یک دقت قابل قبول کوچکتر خواهند بود.

کلمات کلیدی

درخت تصمیم، قوانین دسته بندی، کاهش ویژگی ها

Feature Selection as an Improving Step for Decision Tree Construction

Mehdi Esmaeili; Mansour Tarafdar

Abstract

The removal of irrelevant or redundant attributes could benefit us in making decisions and analyzing data efficiently. Feature Selection is one of the most important and frequently used techniques in data preprocessing for data mining.

In this work, special attention is made on feature selection for classification with labeled data. Here an algorithm is used that arranges attributes based on their importance using two independent criteria. Then, the arranged attributes can be used as input one simple and powerful algorithm for construction decision tree. Results indicate that this decision tree using featured selected by proposed algorithm outperformed decision tree without feature selection. From the experimental results, it is observed that, this method generates smaller tree having an acceptable accuracy.

Key Words

Decision Tree, Classification Rules, Features Reduction

۱. مقدمه

انتخاب صفات خاصه مناسب نقش مهمی در تکنیک های داده کاوی بازی می کند. بدون شک متدها معمولاً با داده های کاهش یافته سریع تر عمل می کنند. وجود داده های تکراری و نامربوط به طور معمول فقط به افزایش حجم داده ها و یا گمراهی متدها منجر می شود.

۱. عضو هیئت علمی دانشگاه آزاد اسلامی کاشان m.esmaeili@iaukashan.ac.ir

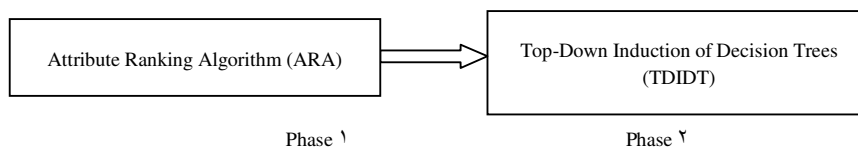
۲. دانشجوی کارشناسی ارشد نرم افزار دانشگاه آزاد اسلامی قزوین، دانشکده برق رایانه و فناوری اطلاعات tarafdar.mansour@gmail.com

از سویی دیگر می‌دانیم که الگوریتم‌های دسته‌بندی یکی از متدهای رایج در داده‌کاوی به حساب می‌آیند. در این بین درخت‌های تصمیم به دلیل سادگی فهم از محبوبیت خوبی برخوردار هستند. هدف اصلی انتخاب صفات خاصه در الگوریتم‌های دسته‌بندی، کاهش صفات خاصه در ساخت درخت تصمیم است به صورتی که یک دقت قابل قبول برای مدل ساخته شده حفظ شود. برای یافتن حالت بهینه از ترکیب صفات خاصه، کلیه حالات آن باید بررسی و تحلیل شوند. روشن است که در بسیاری از مواقع این بررسی بسیار پیچیده بوده و یا حتی شدنی نیست. تحت شرایط خاصی می‌توان از الگوریتم‌هایی بهره برد که نیمه بهینه^۱ هستند. هرچند که از نام این الگوریتم‌ها مشخص است این روش‌ها بهترین گزینه را انتخاب نمی‌کنند اما با انتخاب زیر مجموعه‌ای مناسب از صفات خاصه و پیچیدگی کمتر، کارایی بهتری از خود نشان می‌دهند. جستجو برای انتخاب صفات خاصه می‌تواند به صورت اتفاقی شروع شود (استفاده از الگوریتم‌هایی نظیر الگوریتم ژنتیک)^[۱]. شاخص‌های ارزیابی از نقطه نظر الگوریتم‌هایی که پس از انتخاب صفات خاصه بر روی داده‌ها اجرا می‌شوند به دو دسته مستقل^۲ و غیرمستقل^۳ دسته‌بندی می‌شوند^[۲]. یک معیار مستقل سعی می‌کند تا زیر مجموعه انتخابی از صفات خاصه^۴ را با کمک خصوصیات و خواص خود داده‌ها بدست آورد بدون اینکه به الگوریتم اجرایی در مراحل بعد توجهی داشته باشد. برخی از پرکاربردترین آنها عبارتند از: معیارهای مبتنی بر فاصله^۵، معیارهای اطلاعاتی^۶ و معیارهای سازگاری^{[۲][۳][۴][۵][۶]}. این درحالی است که معیارهای غیرمستقل کاملاً وابسته به الگوریتم اجرایی بعدی هستند. بدین صورت که به منظور ارزیابی روش پس از انتخاب زیر مجموعه‌ای از صفات خاصه، کارایی الگوریتم با کمک این زیرمجموعه انتخابی ارزیابی می‌شود و مناسب‌ترین زیر مجموعه، وابسته به اجرای بهتر الگوریتم با این صفات خاصه خواهد بود. اخیراً تحقیقات زیادی بر روی انتخاب صفات خاصه مناسب جهت دسته‌بندی داده‌ها انجام شده است. لیستی از آنها در [۱][۲][۷][۸] آمده‌اند.

مقاله حاضر از الگوریتمی جهت رتبه‌بندی^۸ صفات خاصه براساس اهمیت‌شان استفاده می‌کند. این الگوریتم از دو معیار مستقل بهره می‌برد. پس از آن صفات خاصه به ترتیب اهمیت‌شان وارد مرحله ساخت درخت تصمیم خواهند شد. این مقاله به صورت زیر سازماندهی شده است: بخش دوم جزئیات الگوریتم پیشنهادی را توضیح می‌دهد. در بخش سوم خصوصیات مجموعه داده‌هایی که برای ارزیابی الگوریتم پیشنهادی انتخاب شده‌اند را به صورت خلاصه آورده‌ایم. بخش چهارم شامل نتایج حاصل از الگوریتم پیشنهادی و مقایسه آن با چند الگوریتم دیگر است و در بخش پایانی نتیجه‌گیری و پیشنهاداتی را برای کارهای آینده بیان کرده‌ایم.

۲. روش پیشنهادی

روش پیشنهادی دارای دو مرحله است که در شکل ۱ نمودار اجمالی آن نشان داده شده است.



شکل ۱: دیاگرام روش پیشنهادی

۲.۱. مرحله اول: رتبه‌بندی صفات خاصه

در این مرحله همانطور که از نام آن مشخص است از الگوریتمی (ARA) جهت بررسی اهمیت صفات خاصه استفاده می‌شود. بدین منظور از معیارهای سنجشی که در [۹] به آن اشاره شده است استفاده کرده‌ایم. در الگوریتم ARA ما به دنبال صفات خاصه‌ای می‌گردیم که با کمک آن می‌توان تا حد امکان دو کلاس را از یکدیگر تمیز داد. همراه با این موضوع از معیار دیگری جهت ارزیابی همبستگی^۹ صفت خاصه به کلاس مربوطه نیز استفاده می‌شود. افزایش و بزرگی در این معیارها اهمیت بیشتر صفت خاصه را نشان می‌دهد به صورتی که می‌توان اهمیت یک صفت خاصه را از برآورد مجموع این دو معیار ارزیابی نمود. ورودی الگوریتم مجموعه داده‌هایی است با تعداد n صفت خاصه در C کلاس قرار گرفته‌اند. با حذف صفت خاصه اول اهمیت این صفت خاصه محاسبه می‌شود. این کار برای n صفت خاصه تکرار می‌شود. عدد بزرگتر نشان دهنده اهمیت بالاتر صفت خاصه است. بدین ترتیب توانسته‌ایم ترتیبی برای صفات خاصه در نظر بگیریم. فرمول‌های ۱ و ۲ نشان می‌دهند که این کار چگونه انجام می‌شود.

$$\text{Distance1} = \sum_{i=1}^c P_i \sum_{k=1}^{n_i} \left[(X_{ik} - m_i)(X_{ik} - m_i)^T \right]^{1/2} \quad (1)$$

$$\text{Distance2} = \sum_{i=1}^c P_i \left[(m_i - m)(m_i - m)^T \right]^{1/2} \quad (2)$$

در این فرمول ها P_i احتمال کلاس i ام را نشان می‌دهد. X ها مقادیر داده ها را نشان می‌دهند که قبل از هر گونه پردازشی نرمال سازی می‌شوند. X_{k_i} مقدار نرمال سازی شده k امین صفت خاصه از i امین نمونه (رکورد) را نشان می‌دهد. m_i و m به ترتیب بردارهایی هستند که میانگین را برای نمونه های کلاس i ام و میانگین کل مجموعه داده ها را مشخص می‌کنند. n_i تعداد نمونه ها در کلاس i است. به صورتی که می‌توان نوشت: $n = n_1 + n_2 + \dots + n_c$ به منظور محاسبه همبستگی صفت خاصه k ام و کلاس ها نیز از فرمول ۳ استفاده می‌کنیم:

$$\text{Attribute class correlation} = \sum_{i \neq j} |X_{ik} - X_{jk}| \quad (3)$$

این فرمول برای صفات خاصه ای که متعلق به یک کلاس نیستند محاسبه می‌شود.

۲, ۲. مرحله دوم: ساخت درخت تصمیم

در این مرحله از یک الگوریتم ساده در عین حال قدرتمند جهت تولید شروط^{۱۰} درخت تصمیم استفاده می‌شود. شکل ۲ این الگوریتم را نشان می‌دهد.

```

IF all the instances in the training set belong to the same class THEN
    Return the value of class
ELSE (a) Select an attribute A from ranked list
      (b) Sort the instances in the training set into subsets,
           one for each value of attribute A
      (c) Return a tree with one branch for each non-empty
           subset, Each branch having a descendant subtree or
           a class value Produced by applying the algorithm
           recursively
  
```

شکل ۲: الگوریتم ساخت درخت تصمیم

۳. توصیف داده های تحت آزمایش

از آنجا که معمولا روش های آماده سازی داده ها وابسته به نوع کاربرد هستند در این قسمت سعی بر آن شده تا داده هایی با خصوصیات متفاوت انتخاب شوند. این مجموعه داده ها از پایگاه داده UCI انتخاب شده اند [۱۰]. در جدول ۱ این چهار مجموعه داده ها به صورت خلاصه توصیف شده اند.

جدول ۱: مشخصات کلی از داده های تحت آزمایش

	of Number Attributes	Number of Instances	Number of Classes
Iris	۴	۱۵۰	۳
Monk's Problems	۷	۴۳۲	۲
Glass Identification	۱۰	۲۱۴	۶
Ionosphere	۳۴	۳۵۱	۲

۴. نتایج تجربی و بحث

در این بخش خروجی الگوریتم پیشنهادی را بر روی چهار مجموعه داده های معرفی شده در بخش قبل بررسی می کنیم. خروجی فاز اول الگوریتم برای مجموعه داده ها در جدول ۲ نشان داده شده است. در این جدول رتبه بندی میان صفات خاصه هر مجموعه داده را نشان می دهد.

جدول ۲: رتبه بندی صفات خاصه (خروجی فاز اول)

Data Set	Attributes Ordering
Iris	۳,۲,۱,۴
Monk's Problems	۲,۵,۴,۱,۳,۶
Glass Identification	۲,۳,۸,۴,۱,۶,۷,۵,۹
Ionosphere	۱۴,۲۰,۲۲,۱۲,۳۰,۴,۲۷,۲۸,۱۸,۱۶,۲۴,۲,۶,۱۰,۲۳,۳۲,۷, ۱۹,۸,۳۱,۲۱,۲۵,۱۷,۱۳,۲۹,۱۱,۲۶,۹,۳,۵,۱۵,۳۳,۳۴,۱

همانطور که مشاهده می کنید در مجموعه داده Iris به ترتیب صفت خاصه سوم و دوم و اول و چهارم اهمیت صفات خاصه را در این مجموعه داده تشکیل می دهند. شکل ۳ خروجی فاز دوم الگوریتم پیشنهادی را بر روی نمونه داده Iris نشان می دهد.

Field^۳ ≤ ۱.۷ : Iris-setosa (۴/۸)
Field^۳ > ۱.۷ & Field^۲ ≤ ۲.۲ : Iris-versicolor (۴/۱)
Field^۳ > ۱.۷ & Field^۲ > ۲.۲ & Field^۱ ≤ ۴.۹ : Iris-versicolor (۳/۲)
Field^۳ > ۱.۷ & Field^۲ > ۲.۲ & Field^۱ > ۴.۹ & Field^۴ ≤ ۱.۴ : Iris-versicolor (۳/۲)
Field^۳ > ۱.۷ & Field^۲ > ۲.۲ & Field^۱ > ۴.۹ & Field^۴ > ۱.۴ : Iris-virginica (۶۱/۱۴)

شکل ۳: خروجی فاز دوم الگوریتم بر روی داده Iris

چنانچه به پنج شرط تولید شده توسط برنامه توجه کنید خواهید یافت که ترتیب آنها همان ترتیب اهمیت صفات خاصه است که در فاز اول الگوریتم به دست آمده است. اعداد داخل پرانتز نشان دهنده تعداد نمونه هایی است که شامل شروط مربوط می شوند. جدول ۳ صحت و دقت روش پیشنهادی بر روی چهار مجموعه داده انتخابی نشان می دهد.

جدول ۳: محاسبه پارامترها به تفکیک کلاس

	TP Rate	FP Rate	Recall	Precision	F-measure	Class
Iris	۰.۹۶	۰	۰.۷۲	۱.۰۰	۰.۸۴	Iris Setosa
	۰.۷۲	۰.۰۵	۰.۷۲	۰.۸۸	۰.۸۰	Iris Versicolour
	۰.۹۴	۰.۱۴	۰.۹۰	۰.۷۷	۰.۸۳	Iris Virginica
Monk's Problems	۰.۶۷	۰.۳۳	۰.۶۷	۰.۶۷	۰.۶۷	Class ۰
	۰.۶۷	۰.۳۳	۰.۶۷	۰.۶۷	۰.۶۷	Class ۱
Glass Identification	۰.۱۰	۰.۰۳	۰.۰۷	۰.۶۴	۰.۱۳	Building_w_f_p
	۰.۹۶	۰.۵۸	۰.۸۴	۰.۴۸	۰.۶۱	Building_w_nf_p
	۰.۰۶	۰.۰۱	۰.۰۱	۰.۳۴	۰.۰۲	Vehicle_w_f_p
	۰.۵۴	۰.۰۱	۰.۰۷	۰.۷۸	۰.۱۳	Containers
	۰.۷۸	۰.۰۲	۰.۰۸	۰.۶۴	۰.۱۴	Tableware
	۰.۸۶	۰.۰۱	۰.۲۲	۰.۹۳	۰.۳۶	Headlamps
Ionosphere	۰.۷۰	۰.۱۰	۰.۷۰	۰.۸۷	۰.۷۸	Bad
	۰.۹۴	۰.۱۷	۰.۹۴	۰.۸۵	۰.۸۹	Good

از ستون آخر جدول ۳ روشن است این مقادیر به صورت مجزا برای هر کلاس محاسبه شده است. بدون شک مدل‌های تولید شده توسط الگوریتم‌های داده کاوی هرچه کوچکتر باشند قابل فهم تر خواهند بود. در دسته بندی داده ها و به خصوص درخت های تصمیم تعداد شروط استخراج شده توسط الگوریتم یکی از مهمترین پارامترهایی به شمار می رود که در مقایسه کارایی الگوریتم ها بر روی آن تاکید می شود. در جدول ۴ مقایسه ای میان چهار الگوریتم و متد پیشنهادی براساس تعداد شروط تولید شده و اندازه درخت انجام شده است.

جدول ۴ : تعداد شروط تولیدشده و اندازه درخت توسط الگوریتمها

	J۴۸	BFTree	REPTree	NBTree	Proposed Method
Iris	۵/۹	۶/۱۱	۳/۵	۴/۷	۵/۸
Monk's Problems	۲/۳	۲/۳	۸/۱۵	۱/۱	۴/۶
Glass Identification	۳۰/۵۹	۱۶/۳۱	۱۲/۲۳	۹/۱۷	۱۴/۳۰
Ionosphere	۱۸/۳۵	۱۱/۲۱	۵/۹	۸/۱۵	۱۳/۲۴

همانطور که از جدول مشخص است روش پیشنهادی در بین این روشها بهترین نیست اما عکس العمل آن نسبت به نوع داده ها خوب است. توجه کنید که ما در ساخت درخت از هیچ راهبرد اکتشافی^{۱۱} استفاده نمی کنیم و این درحالی است که دیگر الگوریتم ها پیچیده تر از فاز دوم الگوریتم ما عمل می کنند. در جدول ۵ نرخ خطای روش پیشنهادی نیز بررسی شده است.

جدول ۵ : مقایسه نرخ خطای الگوریتم ها

	Error Rate				
	J۴۸	BFTree	REPTree	NBTree	Proposed Method
Iris	۰.۰۴	۰.۰۶	۰.۰۶	۰.۰۶	۰.۱۲
Monk's Problems	۰.۲۵	۰.۲۵	۰.۱۵	۰.۲۵	۰.۳۳
Glass Identification	۰.۳۴	۰.۳۳	۰.۳۸	۰.۳۰	۰.۴۵
Ionosphere	۰.۰۹	۰.۱۰	۰.۱۱	۰.۱۰	۰.۱۸

۵. نتیجه گیری

در این مقاله تاثیر مرحله پیش پردازش و آماده سازی داده ها در ساخت یک درخت تصمیم بررسی شد. انتخاب صفات خاصه مناسب برای ساخت یک درخت تصمیم نقش مهمی بازی می کند. در الگوریتم پیشنهادی ابتدا ترتیبی میان صفات خاصه در داده ها به دست می آید. همبستگی داده ها به برچسب کلاس ها و تغییر مقادیر داده ها با تغییر کلاس از جمله پارامترهایی است که در اهمیت رتبه بندی یک صفت خاصه محاسبه می شود. در مرحله بعدی با توجه به این رتبه بندی درخت تصمیم و یا در واقع شروط ایجاد می شوند. نتایج به دست آمده نشان می دهد درختی که توسط روش پیشنهادی تولید می شود به طور متوسط کوچکتر از درخت هایی است که توسط چهار متد انتخابی ایجاد شده است. چنانچه مایل باشیم تا درخت حاصل از الگوریتم پیشنهادی دارای دقت بالاتر و بهتری باشد می توانیم مرحله اول الگوریتم را پس از حذف صفت خاصه به صورت تکراری انجام دهیم. در این صورت الگوریتم از پیچیدگی بیشتری برخوردار و در نتیجه سرعت اجرای آن کاهش می یابد. بهر حال به عنوان کار آینده می توان افزودن یک روش هرس کردن مناسب برای درخت را آزمایش نمود. بدون شک چون درخت تولید شده توسط روش پیشنهادی ترتیبی در تست کردن شروط دارد (Oblivious Tree) شاید بتوان قوانین مناسب تر و بهتری برای هرس کردن آن با حفظ دقت مناسب جستجو نمود.

- [۱] J.Doak , *An Evaluation of Feature Selection Methods and Their Application to Computer Security* , technical report, University of California at Davis , Department of Computer Science, ۱۹۹۲
- [۲] H.Liu, L.Yu , *Toward Integrating Feature Selection Algorithms for Classification and Clustering* , IEEE Transactions on Knowledge and Data Engineering , Vol. ۱۷, No. ۴, pp. ۴۹۱-۵۰۲, April-۲۰۰۵
- [۳] H.Almuallim , T.G. Dietterich , *Learning Boolean Concepts in the Presence of Many Irrelevant Features* , Artificial Intelligence, Vol. ۶۹, pp.۲۷۹-۳۰۵, ۱۹۹۴
- [۴] M.Ben-Bassat, *Pattern Recognition and Reduction of Dimensionality* , Handbook of Statistics-II, P.R. Krishnaiah and L.N. Kanal, pp. ۷۷۳-۷۹۱, North Holland, ۱۹۸۲
- [۵] M.A.Hall , *Correlation-Based Feature Selection for Discrete and Numeric Class Machine Learning* , Proc. ۱۷th Int'l conf. Machine Learning, pp. ۳۵۹-۳۶۶, ۲۰۰۰
- [۶] H.Liu , H.Motoda , *Feature Selection for Knowledge Discovery and Data Mining. Boston* , Kluwer Academic, ۱۹۹۸.
- [۷] M.Dash , H.Liu , *Feature Selection for classification* , Intelligent Data Analysis: An Int'l J., Vol. ۱, No. ۳, pp. ۱۳۱-۱۵۶, ۱۹۹۷
- [۸] W.Siedlecki , J.Sklansky , *On Automatic Feature Selection* , Int'l J. Pattern Recognition and Artificial Intelligence, Vol. ۲, pp. ۱۹۷-۲۲۰, ۱۹۸۸.
- [۹] L.Wang, X.Fu, *Data Mining with Computational Intelligence* , Springer , pp. ۱۱۷-۱۲۳, ۲۰۰۵ .
- [۱۰] C.L. Blake, C.J.Merz, *UCI Repository of Machine Learning Databases*. Irvine, CA: University of California, Department of Information and Computer Science , (۱۹۹۸). [<http://www.ics.uci.edu/~mlearn/MLRepository.html>]

^۱ Suboptimal

^۲ Independent

^۳ Dependent

^۴ Feature Subset

^۵ Distance Measures

^۶ Information Measures

^۷ Consistency Measures

^۸ Ranking

^۹ Correlation

^{۱۰} Rules Generation

^{۱۱} Heuristic approach