



تعیین مهم ترین ویژگی های کمی و کیفی متمایز کننده سرده *Rubus L.* در ایران با استفاده از الگوریتم های دسته بندی و انتخاب ویژگی

محمدجواد شیخزاده*

دانشکده علوم پایه و فنی مهندسی، دانشگاه گنبد کاووس، گنبد کاووس، ایران

* Email: sheikhzadeh@gonbad.ac.ir

تاریخ پذیرش: ۹۷/۰۴/۲۰

تاریخ دریافت: ۹۷/۰۱/۰۷

چکیده

جنس سرده *Rubus L.* متعلق به تیره Rosaceae و زیر تیره Rosoideae شامل حدود ۷۵۰ گونه در دنیا است. این جنس در اکثر نواحی جهان پراکنش دارد. در فلور ایران تعداد هشت گونه و پنج هیبرید (دو رگه) گزارش شده است. تمشک یکی از پر چالش ترین جنس ها در بین گیاهان گل دار می باشد. وجود عواملی از جمله پلی پلوئیدی، آپومیکیسی و دورگه گیری زیاد در این جنس باعث مشکلاتی در تشخیص از نظر ریخت شناسی شده است. جمع آوری داده های کمی و کیفی جهت تشخیص گونه ها و روش های جمع آوری مولفه های ریخت شناسی بسیار زمانبر و پر هزینه است. بنابراین، بکارگیری روش های متفاوت در جهت کاهش زمان و هزینه ها همواره مطرح می باشد. در بسیاری از موارد، جهت آنالیز داده های زیستی روش های داده کاوی بکار گرفته می شود. در این مقاله، از ترکیب الگوریتم های مختلف انتخاب ویژگی و دسته بندی برای تشخیص ویژگی های متمایز کننده بین گونه های سرده *Rubus L.* استفاده شد. با بکارگیری روش دسته بندی Random Forest و مدل انتخاب ویژگی InfoGainAttributeEval با تعداد ۲۸ ویژگی به دقت ۹۴/۰۵ درصد در دسته بندی رسیدیم که بهترین الگوریتم از نظر دقت می باشد و با استفاده از روش MLP و مدل انتخاب ویژگی SymmetricalAttributeEval، با تعداد ۴ ویژگی دقت دسته بندی ۸۴/۳۲ درصد حاصل شد که بهترین الگوریتم از نظر تعداد اندک ویژگی های انتخاب شده است. چهار ویژگی فوق توسط اکثر الگوریتم های استفاده شده در این مقاله انتخاب شدند. تمامی این ویژگی ها کیفی هستند و جهت بدست آوردن آنها نیازی به هزینه اندازه گیری آزمایشگاهی نمی باشد. بنابراین می توانند معیار مناسبی جهت کلید شناسایی باشند.

کلیدواژه ها: الگوریتم، داده کاوی، ریخت شناسی، کلید شناسایی.

مقدمه

است [۱۲،۲۲،۲۹،۳۱] که در اکثر نواحی جهان به جز ناحیه قطب جنوب پراکنش دارد [۹،۱۰،۱۱،۱۶]. این تاکسون از جنبه اقتصادی و بوم شناسی حائز اهمیت

جنس تمشک متعلق به تیره Rosaceae و زیر تیره Rosoideae می باشد و شامل حدود ۷۵۰ گونه در دنیا

زمان صرف شده برای استخراج و همچنین ارزیابی آنها کاهش می‌یابد. وجود داده‌های بزرگ در برخی از مسائل که ارزیابی و تحلیل آنها سخت، پرهزینه و زمان‌بر است و دقت پایین الگوریتم‌ها باعث می‌شود تا همچنان پیدا کردن الگوریتمی که بتواند به طوری کارآمد ویژگی‌های تفکیک‌کننده را مشخص کرده و با استفاده از آنها نمونه‌ها را دسته‌بندی کند حائز اهمیت باشد [۳۳،۳۴].

هدف از این پژوهش، معرفی روشی بهینه برای تشخیص مهم‌ترین ویژگی‌های متمایزکننده و همچنین دسته‌بندی گونه‌های مختلف سرده تمشک در ایران است. طوری که با استفاده از تعدادی کمتر از ویژگی‌ها با دقت بیشتری نمونه‌ها را دسته‌بندی کرد. در بسیاری از موارد، جهت آنالیز داده‌های زیستی (از جمله مطالعات آناتومی و ریخت‌شناسی) روش‌های داده کاوی بکار گرفته می‌شود. بنابراین در این مقاله، از ترکیب الگوریتم‌های مختلف انتخاب ویژگی و دسته‌بندی برای تشخیص ویژگی‌های متمایزکننده بین گونه‌های سرده *Rubus L.* استفاده می‌شود.

مواد و روش‌ها

در این مقاله، به منظور انجام مطالعات ریخت‌شناسی روی جمعیت‌های مختلف گونه‌های معرفی شده سرده *Rubus* در ایران (شامل *R. caesius*، *R. hirtus*، *R. dolichocarpus*، *R. discolor*، *R. persicus* و *R. sanctus*) از مناطق مختلف شمال ایران از سه استان گیلان، مازندران و گلستان نمونه‌هایی مطابق (جدول ۱) جمع‌آوری گردید. به این ترتیب در این پژوهش از ۵۷۸۱ داده خام حاصل از تعداد ۴۶ ویژگی کیفی و کمی استفاده شده است، که این ویژگی‌ها در (جدول ۲ و ۳) ارائه

است، چون دارای میوه‌های خوراکی و گیاهانی با ارزش تزئینی هستند. این جنس به واسطه داشتن فرم رویشی درختچه‌ای و رونده، ساقه خاردار و میوه خوراکی از نوع شفت مجتمع متشکل از شفتچه‌هایی کوچک به آسانی قابل تشخیص است. اگرچه تشخیص گونه‌های آن به واسطه تنوع ریخت‌شناسی در جمعیت‌های گونه بسیار مشکل است. این جنس در ۱۲ زیرجنس طبقه‌بندی می‌شود [۹-۱۱]. بزرگترین زیرجنس آن *Rubus* با ۶ بخش و ۱۳۲ گونه می‌باشد و کوچکترین زیر جنس آن *Chamaemorus* با یک گونه می‌باشد. زیر جنس *Rubus* شامل گونه‌هایی هستند که در اروپا، آسیا و آمریکای شمالی یافت می‌شوند [۳]. خاتم‌ساز هشت گونه و پنج هیبرید (دورگه) را از ایران نام برده است [۲۰]، که به دو زیر جنس *Cylactis* و *Rubus* تعلق دارد. *R. saxatilis* به‌عنوان تنها گونه‌ی علفی جنس در ایران متعلق به زیرجنس *Cylactis* و ۷ گونه *R. discolor*، *R. caesius*، *R. dolichocarpus*، *R. hirtus*، *R. hyrcanus*، *R. persicus* و *R. sanctus* متعلق به زیرجنس *Rubus* می‌باشند. تمشک یکی از پرچالش‌ترین جنس‌ها در بین گیاهان گل‌دار است [۱،۲۲،۲۷،۲۹] که این چالش‌ها به خاطر وجود عوامل بیولوژیکی از جمله پلی‌پلوئیدی، آپومیکیسی و دورگه‌گیری زیاد بین گونه‌های یک زیر-جنس یا گونه‌های زیرجنس‌های مختلف بوده و به نظر می‌رسد استفاده از ریخت‌شناسی ظاهری به‌تنهایی در روشن شدن وضعیت رده‌بندی این جنس کارایی نداشته باشد [۱۳،۳۲]. از سوی دیگر، استخراج ویژگی‌های کمی و کیفی یک گیاه کاری بسیار زمان‌بر و پرهزینه است. از طرفی، بسیاری از آنها یا بلااستفاده هستند و یا اینکه بار اطلاعاتی چندانی ندارند. از این رو هرچه تعداد این ویژگی‌ها کمتر باشد، هزینه و

گردیده است. ویژگی‌ها کدگذاری شده و بر اساس است [۱۸].
تعداد نمونه‌ها، میانگین و واریانس محاسبه گردیده

جدول ۱- نمونه‌های جمع‌آوری شده در سه استان مورد مطالعه [۱۸]

ردیف	گونه	محل جمع‌آوری	ارتفاع (متر)	جمع‌آوری کننده	مشخصات هرباریومی
۱	<i>R. caesius</i> L	مازندران، رامسر، قاسم آباد سفلی	۵۰	کسلخه، حبیبی	-803012 GKUH
۲	<i>R. caesius</i> L	گلستان، بندر گز، گز غربی	۹۵	کسلخه، حبیبی	-803025 GKUH
۳	<i>R. discolor</i> Weihe and Nees.	گیلان، جاده آستارا به اردبیل، گردنه حیران	۵۰۰	کسلخه، مهدی‌یانی	-803055 GKUH
۴	<i>R. dolichocarpus</i> Jaz.	گلستان، پارک ملی گلستان	۵۰۰	کسلخه، مهدی‌یانی	803120- GKUH
۵	<i>R. hirtus</i> Waldst and kit	مازندران، جنگل سنگده	۱۳۶۰	کسلخه، مهدی‌یانی	-803135 GKUH
۶	<i>R. hyrcanus</i> Juz.	گیلان، جاده اسالم به خلخال	۱۰۳۰	کسلخه، حبیبی	-803146 GKUH
۷	<i>R. persicus</i> Boiss	گلستان، جنگل امام رضا (ع) کردکوی	۲۵۰	مهدی‌یانی	-803150 GKUH
۸	<i>R. sanctus</i> Schreber	گیلان، جیرنده، روستای بی‌ورزین، نزدیک زیارتگاه	۱۰۴۰	کسلخه، حبیبی	-803235 GKUH

جدول ۲- لیست ویژگی‌های کیفی سرده *Rubus* و کد گذاری آنها در مطالعات ریخت‌شناسی [۱۸]

ردیف	ویژگی‌ها	شرح	انحراف معیار	واریانس
۱	فرم رویشی	بوته‌ای (۰)، درختچه‌ای (۱)	۰/۳۷۵۰	۰/۸۳۰۶
۲	فرم ساقه	خوابیده (۰)، خمیده (۱)، افراشته (۲)، بالارونده (۳)	۰/۷۶۴۵	۱/۴۳۵۴
۳	شکل شاخه گل‌دهنده	گرد (۰)، تا حدودی زاویه‌دار (۱)، زاویه‌دار (۲)	۰/۷۶۳۸	۱/۴۲۷۴
۴	شکل خار شاخه گل‌دهنده	مستقیم کوتاه (۰)، خمیده کوچک (۱)، خمیده درشت (۲)، سوزنی (۳)، خمیده سوزنی (۴)، خمیده و مستقیم (۵)	۱/۳۰۲۰	۲/۲۶۶۱
۵	حضور کرک در شاخه گل‌دهنده	بدون کرک (۰)، تا حدودی کرک‌دار (۱)، کرک‌دار (۲)	۰/۳۵۲۲	۱/۸۵۴۸
۶	نوع کرک در شاخه گل‌دهنده	تار (۰)، ستاره (۱)، تار و ستاره (۲)، نمدی (۳)	۰/۵۴۶۲	۲/۵
۷	حضور غده پایکدار در شاخه	حضور غده پایکدار (۰)، عدم حضور غده (۱)، حضور غده بدون پایک (۲)	۰/۳۷۴۴	۰/۸۵۴۸
۸	شکل برگچه جانبی	تخم‌مرغی (۰)، واژتخم‌مرغی (۱)، بیضوی (۲)، بیضوی پهن (۳)، لوزی (۴)، دایره‌ای (۵)، تخم‌مرغی - بیضوی (۶)، واژتخم‌مرغی - بیضوی (۷)	۲/۷۹۳۹	۳/۹۸۳۸
۹	شکل برگچه انتهایی	تخم‌مرغی (۰)، واژتخم‌مرغی (۱)، بیضوی (۲)، لوزی (۳)، بیضوی پهن با عرض بیش از ۲.۵ برابر طول (۴)، دایره‌ای (۵)	۱/۲۶۵۴	۱/۷۸۶۲
۱۰	نوک برگچه جانبی	بدون نوک (۰)، تیز یا کمی کشیده (۱)، دم‌دار (۲)، نوک کند (۳)	۰/۸۳۷۳	۱/۳۷۰۹
۱۱	نوک برگچه انتهایی	بدون نوک (۰)، تیز یا کمی کشیده (۱)، دم‌دار (۲)، نوک کند (۳)	۰/۸۴۴۰	۱/۴۲۷۴

ردیف	ویژگی‌ها	شرح	انحراف معیار	واریانس
۱۲	قاعده برگچه جانبی	قلبی (۰)، دایره‌ای (۱)، قلبی - دایره‌ای (۲)	۰/۶۱۲۲	۱/۵۶۴۵
۱۳	قاعده برگچه انتهایی	قلبی (۰)، دایره‌ای (۱)، قلبی - دایره‌ای (۲)	۰/۸۱۶۱	۱/۵۵۶۴
۱۴	رنگ سطح تحتانی برگچه	سبز (۰)، سفید (۱)، تا حدودی سبز (۲)، تا حدودی سفید (۳)	۰/۶۷۶۷	۰/۹۵۹۶
۱۵	پوشش کرک سطح فوقانی برگچه	بدون کرک (۰)، پراکنده (۱)، متوسط (۲)، فشرده (۳)	۰/۷۵۸۵	۱/۱۹۳۵
۱۶	پوشش کرک سطح تحتانی برگچه	بدون کرک (۰)، پراکنده (۱)، متوسط (۲)، فشرده (۳)	۰/۵۹۳۲	۲/۶۰۴۸
۱۷	نوع کرک سطح فوقانی برگچه	بدون کرک (۰)، اکثرا تار (۱)، اکثرا ستاره (۲)، تار و ستاره (۳)، نمدی (۴)	۱/۰۷۳۳	۱/۴۶۷۷
۱۸	نوع کرک سطح تحتانی برگچه	بدون کرک (۰)، اکثرا تار (۱)، اکثرا ستاره (۲)، نمدی (۳)	۰/۰۸۹۴	۲/۹۹۱۹
۱۹	شکل خار دمبرگ	مستقیم کوتاه (۰)، خمیده کوچک (۱)، خمیده درشت (۲)، سوزنی (۳)، خمیده سوزنی (۴)	۰/۹۵۳۴	۱/۱۸۵۴
۲۰	نوع کرک دمبرگ	اکثرا تار (۰)، اکثرا ستاره (۱)، نمدی (۳)، تار و ستاره (۴)	۰/۵۴۱۴	۲/۴۲۷۴
۲۱	وضعیت دمبرگچه جانبی	بدون دمبرگچه (۰)، دارای دمبرگچه با طول بیش از ۲mm (۱)، دمبرگچه خیلی کوتاه با طول کمتر از ۲mm (۲)	۰/۹۴۶۹	۱/۱۲۰۹
۲۲	نوع گل‌آذین	دیهم (۰)، خوشه ساده (۱)، خوشه گرزن مرکب (۲)، خوشه مرکب (۳)	۰/۹۳۹۳	۱/۶۱۲۹
۲۳	حضور غده در گل‌آذین	غده‌دار (۰)، بدون غده (۱)	۰/۴۰۷۰	۰/۷۹۰۳
۲۴	شکل کاسبرگ	تخم‌مرغی (۰)، مستطیلی (۱)	۰/۹۸۱۲	۰/۸۵۴۸
۲۵	نوک کاسبرگ	بدون نوک تا کمی (۰)، تیز (۱)، سیخکی (۲)	۰/۹۵۰۷	۰/۸۲۲۵
۲۶	نوع کرک کاسبرگ	نمدی (۰)، تار و ستاره (۱)	۰/۴۲۳۲	۰/۲۳۳۸
۲۷	غده پایک‌دار در کاسبرگ	غده‌دار (۰)، بدون غده (۱)	۰/۴۲۳۲	۰/۷۶۶۱
۲۸	شکل گلبرگ	بیضوی کشیده (۰)، بیضوی واژگون (۱)	۱/۰۷۳۱	۰/۵۴۰۳
۲۹	رنگ گلبرگ	صورتی (۰)، سفید (۱)، سفید - صورتی (۲)	۰/۵۱۳۶	۰/۴۵۱۶
۳۰	رنگ میله	صورتی (۰)، سفید (۱)	۰/۴۷۳۲	۰/۶۶۱۲
۳۱	کرک بساک	بدون کرک (۰)، کرک‌دار (۱)	۰/۴۷۵۹	۰/۳۴۶۷
۳۲	شکل گوشوارک	خطی (۰)، خنجری - نیزه‌ای با عرض کمتر از ۲mm (۱)، نیزه‌ای پهن با عرض بیش از ۲mm (۲)	۰/۶۹۷۸	۰/۶۴۵۱
۳۳	تعداد برگچه	یک (۰)، دو الی بیست (۱)، بیست الی چهل (۲)، چهل الی شصت (۳)	۰/۷۴۹۹	۱/۷۴۱۹

جدول ۳- لیست ویژگی‌های کمی سرده *Rubus* و کد گذاری آنها در مطالعه ریخت‌شناسی [۱۸]

ردیف	ویژگی‌ها	انحراف معیار	واریانس
۱	طول دمبرگ	۱/۷۸۴۴	۴/۱۲۴۶
۲	طول برگچه انتهایی	۲/۱۵۴۴	۵/۹۸۵
۳	طول برگچه جانبی	۲/۰۰۹	۵/۱۲۹۳
۴	عرض برگچه انتهایی	۱/۳۷۴۱	۴/۰۹۲۱
۵	عرض برگچه جانبی	۱/۰۵۴۵	۳/۱۵۰۵
۶	طول دمبرگچه برگچه انتهایی	۹/۳۵۴۷	۱۴/۶۲۵۹
۷	طول دمبرگچه برگچه جانبی	۲/۳۸۶۱	۲/۵۱۹۵
۸	طول کاسبرگ	۱/۱۴۹۵	۶/۱۸۹۰

ردیف	ویژگی‌ها	انحراف معیار	واریانس
۹	نسبت طول به عرض کاسبرگ	۰/۴۲۲۳	۲/۱۲۲۰
۱۰	نسبت طول به عرض گلبرگ	۷/۶۳۲۵	۱۲/۳۲۴۵
۱۱	تعداد گل در گل آذین	۱۱/۱۸۳۹	۱۹/۵۰۳۰
۱۲	طول گل آذین	۵/۷۲۰۶	۱۱/۶۴۹۲
۱۳	طول دمگل	۱۰/۴۵۰۷	۱۸/۲۶۶۸

انتخاب شده، از الگوریتم‌های دسته‌بندی استفاده می‌شود.

الگوریتم‌های انتخاب ویژگی^۱

در عصر ارتباطات و اطلاعات، جمع آوری و ذخیره داده‌ها بسیار آسان و ارزان می‌باشد. از سوی دیگر، با افزایش اطلاعات قابل خواندن توسط ماشین، امکان درک و استفاده از اطلاعات بطور همزمان بسیار مشکل بود. ماشین یادگیری ابزارهایی را فراهم می‌کند که می‌توان داده‌های بزرگ را بصورت خودکار تجزیه و تحلیل نمود. لذا اساس یادگیری ماشین، انتخاب ویژگی است. در انتخاب ویژگی باید از بین مجموعه ویژگی‌های مناسب با استفاده از یک روش تفکیک کننده، برجسته ترین و بهترین ویژگی‌ها برای یادگیری انتخاب می‌شود. در روش‌های مختلف انتخاب ویژگی، زیرمجموعه‌ای با کمترین تعداد ویژگی ممکن (ویژگی-های تاثیرگذار) هدف مورد نظر را پیدا می‌کند. فرآیند انتخاب ویژگی در تمامی روش‌ها به بخش‌های زیر تقسیم بندی می‌شود (شکل ۱):

- ۱- روش جستجو (Search Method): تابع تولید کننده زیرمجموعه‌های کاندید شده را برای روش مورد نظر پیدا می‌کند.
- ۲- مدل ارزیابی (Attribute Evaluator): زیر مجموعه کاندید شده را بر اساس روش جستجو ارزیابی

همانگونه که در مقدمه اشاره شد، مساله انتخاب ویژگی، در بسیاری از کاربردها مانند دسته‌بندی (Classification) اهمیت بسیار زیادی دارد، زیرا در این موارد تعداد زیادی ویژگی وجود دارد، که بسیاری از آنها بار اطلاعاتی چندانی ندارند و یا حتی دقت الگوریتم دسته‌بندی را کاهش می‌دهند. بنابراین حذف ویژگی‌های اضافی نه تنها هزینه جمع‌آوری داده‌های کیفی و کمی را کاهش می‌دهد و بار محاسباتی الگوریتم دسته‌بندی را پایین می‌آورد، بلکه ممکن است موجب افزایش دقت دسته‌بندی شود. جهت کاهش ویژگی‌های اضافی و دسته‌بندی آن‌ها، نرم‌افزاری به زبان جاوا پیاده‌سازی شد که از کتابخانه Weka (نسخه ۳/۸) بهره می‌برد که دارای امکانات بسیار گسترده جهت داده‌کاوی (همانند انواع الگوریتم‌های رگرسیون، کاوش قواعد انجمنی، انتخاب ویژگی، خوشه‌بندی و دسته‌بندی داده‌ها)، امکان مقایسه خروجی روشهای مختلف با هم، راهنمای خوب، واسط گرافیکی کارا و سازگاری با سایر برنامه‌ها است [۱۷]. این سیستم به زبان جاوا نوشته شده و بر اساس لیسانس عمومی و فراگیر در پایگاه [۲۳] انتشار یافته است. در نرم‌افزار پیاده‌سازی شده، داده‌ها با استفاده از تعدادی از الگوریتم‌های انتخاب ویژگی واکاوی شده و هر الگوریتم تعدادی ویژگی را از بین ۴۶ ویژگی اولیه انتخاب می‌کند. سپس برای ارزیابی کارایی ویژگی‌های

^۱ Feature Selection Algorithms

ویژگی نیست، اما در عمل باید یک زیرمجموعه ویژگی را در شرایط مختلف امتحان نمود تا شرایط مورد نیاز برای حل مساله مورد نظر احراز شود [۷].

در بخش انتخاب ویژگی کتابخانه weka سه نوع روش جستجو وجود دارد که شامل Bestfirst، GeridyStepwise و Ranker می‌باشد. این کتابخانه شامل مدل‌های ارزیابی مختلف است که در این تحقیق مطابق (شکل ۱) از مدل‌های CfsSubsetEval، InfoGainAttributeEval، CorrelationAttributeEval، SymetricalAttributeEval و ReliefFAttributeEval استفاده شده است [۲۴، ۲۳].

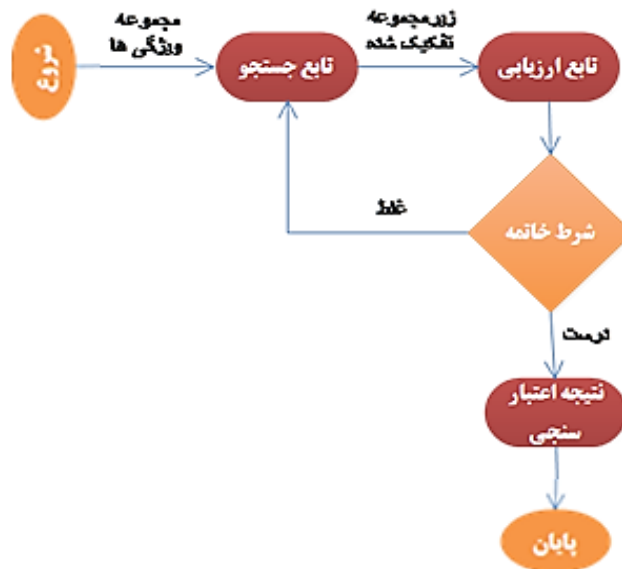
کرده و یک عدد به عنوان میزان مزیت روش باز می‌گرداند. روش‌های مختلف تلاش می‌کنند تا زیرمجموعه‌ای را پیدا کنند که این مقدار را بهینه کند.

۳- شرط خاتمه: برای تصمیم‌گیری در مورد زمان توقف الگوریتم استفاده می‌شود (شکل ۲).

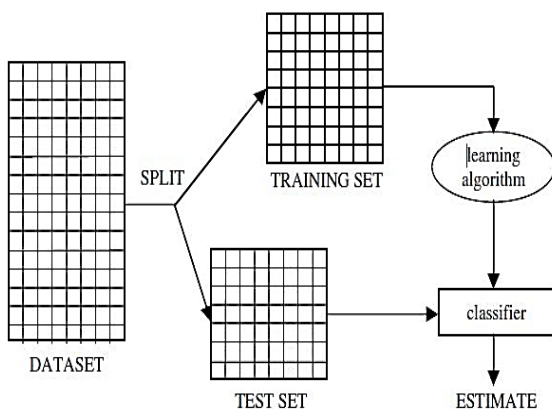
سپس به کمک تابع تعیین اعتبار، اعتبار زیرمجموعه انتخاب شده را تعیین می‌کند. برای این کار می‌توان از داده‌های نمونه برداری شده و یا مجموعه داده‌های شبیه‌سازی شده استفاده نمود. به بیان دیگر، تابع تعیین اعتبار جزئی از فرآیند انتخاب



شکل ۱- فرآیند گزینش بهترین الگوریتم انتخاب ویژگی.



شکل ۲- نمودار فعالیت الگوریتم انتخاب ویژگی در سرده Rubus.



شکل ۳- عملکرد الگوریتم‌های دسته‌بندی [۱۹، ۴].

نتایج و بحث

در این پژوهش نرم‌افزاری با استفاده از کتابخانه Weka و به زبان جاوا نوشته شد که داده‌ها را از فایل ورودی دریافت کرده و عملیات پردازشی (شامل انتخاب ویژگی و دسته‌بندی) را روی داده‌ها اعمال می‌کند، سپس گزارشی (شامل ویژگی‌های انتخاب شده و دقت و خطای دسته‌بندی با استفاده از آن ویژگی‌ها) را به صورت یک فایل اکسل بر می‌گرداند. جهت کاهش تعداد ویژگی‌ها، داده‌ها با استفاده از تعدادی الگوریتم‌های انتخاب ویژگی واکاوی شده و

الگوریتم‌های دسته‌بندی^۱

یکی از مهمترین روش‌های ارزیابی کارایی الگوریتم‌های انتخاب ویژگی استفاده از الگوریتم‌های دسته‌بندی است، به طوری که هرچه دقت دسته‌بندی با استفاده از مجموعه‌ای از ویژگی‌ها بیشتر باشد بیانگر این است که آن ویژگی‌ها از کارایی بیشتری برخوردارند.

جهت ساخت یک دسته‌بندی از روش آموزش و تست (Train and Test Method) مطابق (شکل ۳) استفاده می‌شود که به موجب آن مجموعه داده‌های در دسترس به مجموعه‌های آموزشی (Training Set) و مجموعه‌های تست (Testing Set) تقسیم می‌شود [۱۹، ۴]. در این تحقیق مطابق (شکل ۱) از الگوریتم‌های دسته‌بندی ^۲MLP، Naive Bayes، ^۳C4.5، ^۴KNN و Random Forest استفاده شده است [۱۴، ۲۱، ۲۵، ۲۶، ۲۸، ۳۰].

^۱ Classification Algorithms

^۲ Multilayer Perceptron

^۳ Sequential Minimal Optimization

^۴ K-Nearest Neighbors

هر الگوریتم تعدادی ویژگی را از بین ۴۶ ویژگی اولیه انتخاب نمود. روش جستجوی Ranker به هر ویژگی یک امتیاز اختصاص می‌دهد که هرچه رتبه یک ویژگی بیشتر باشد بدین معنی است که آن ویژگی از اهمیت بیشتری برخوردار می‌باشد. بنابراین جهت انتخاب ویژگی‌های برتر، یک حد آستانه (Threshold) برای رتبه تعیین شده است. به این صورت که اگر امتیاز یک ویژگی از حد آستانه بیشتر شد، آن ویژگی انتخاب می‌شود. حد آستانه نیز بدین شکل بدست آمده است که ویژگی‌های انتخاب شده توسط آن حد آستانه بیشترین دقت را در دسته‌بندی داشته باشند.

تمامی ویژگی‌های انتخاب شده در تجزیه واریانس در سطح یک درصد معنی‌دار می‌باشند. در الگوریتم‌های دسته‌بندی، از ۷۰ درصد داده‌ها (۸۶ نمونه) جهت آموزش و ۳۰ درصد (۳۷ نمونه) باقی مانده برای حالت آزمایش استفاده شده است. به منظور افزایش دقت اعتبارسنجی نتایج، روش Repeated

Random Subsampling با ۲۰۰ تکرار به کار گرفته شده است [۱۵]. این روش اعتبارسنجی استفاده شده است زیرا تمامی حالت‌های ممکن برای آزمایش و آموزش الگوریتم دسته‌بندی را در نظر می‌گیرد. (جدول ۴) میزان دقت و خطای الگوریتم‌های مختلف را نشان می‌دهد. برای ارزیابی دسته‌بندی‌های ذکر شده و انتخاب بهترین دسته‌بندی، از ضریب کاپای کوهن [۵] و ریشه خطای مربع نسبی [۸] استفاده می‌شود. ضریب کاپای کوهن یک معیار آماری توافق درون ارزیاب بین دو اندازه‌گیری برای دسته‌بندی موارد مشابه است. بنابراین هرچه ضریب کاپای کوهن یک الگوریتم بیشتر باشد، میزان دقت آن الگوریتم بالاتر است. همچنین، معیار ریشه خطای مربع نسبی، خطا بر حسب درصدی از میانگین واقعی دقت یک الگوریتم است که هرچه این معیار کمتر باشد میزان دقت آن الگوریتم بیشتر است [۸، ۶].

جدول ۴- نتیجه دسته‌بندی الگوریتم‌ها

الگوریتم دسته‌بندی	الگوریتم انتخاب ویژگی		کاپای کوهن	ریشه مربع خطای نسبی (%)	تعداد ویژگی‌های انتخاب شده	دقت دسته‌بندی	
	روش جستجو	مدل ارزیابی				درصد درست	درصد غلط
C4.5	Best-First	Cfs-Subset-Eval	۰/۸۳۳	۵۴/۵۷	۱۷	۸۷/۲۹	۱۲/۷۱
	Greedy-Stepwise	Cfs-Subset-Eval	۰/۸۵۴	۵۱/۱۴	۲۱	۸۸/۹۱	۱۱/۰۹
	Ranker	Info-Gain-Attribute-Eval	۰/۸۱۸	۵۴/۹۹	۶	۸۵/۹۴	۱۴/۰۶
		Gain-Ratio-Attribute-Eval	۰/۸۳۶	۵۴/۱۸	۸	۸۷/۵۶	۱۲/۴۴
		ReliefF-Attribute-Eval	۰/۷۶۷	۶۱/۵۸	۷	۸۲/۴۳	۱۷/۵۷
		Correlation-Attribute-Eval	۰/۷۶۴	۵۷/۹۲	۵	۸۱/۸۹	۱۸/۱۱
		Symmetrical-Attribute-Eval	۰/۷۵۸	۶۳/۱۶	۱۲	۸۱/۳۵	۱۸/۶۵
	-	-	۰/۸۰۳	۵۸/۲۸	۴۶	۸۴/۵۹	۱۵/۴۱
SMO	Best-First	Cfs-Subset-Eval	۰/۸۷۵	۹۲/۰۵	۱۷	۹۰/۵۴	۹/۴۶
	Greedy-Stepwise	Cfs-Subset-Eval	۰/۸۸۵	۹۱/۹۶	۲۱	۹۱/۳۵	۸/۶۵
	Ranker	Info-Gain-Attribute-Eval	۰/۸۶۷	۹۲/۰۴	۲۳	۸۹/۹۹	۱۰/۰۱

		Gain-Ratio-Attribute-Eval	۰/۸۹۶	۹۱/۹۱	۲۹	۹۲/۱۶	۷/۸۴
		ReliefF-Attribute-Eval	۰/۸۴۳	۹۲/۲۹	۱۸	۸۸/۱۰	۱۱/۹
		Correlation-Attribute-Eval	۰/۹۰۰	۹۱/۹۴	۴۳	۹۲/۴۳	۷/۵۷
		Symmetrical-Attribute-Eval	۰/۸۸۵	۹۱/۹۵	۲۶	۹۱/۳۵	۸/۶۵
	-	-	۰/۸۹۵	۹۱/۲۸	۴۶	۹۱/۸۹	۸/۱۱
Nave Bayes	Best-First	Cfs-Subset-Eval	۰/۹۰۷	۴۱/۰۹	۱۷	۹۲/۹۷	۷/۰۳
	Greedy-Stepwise	Cfs-Subset-Eval	۰/۹۰۰	۴۲/۳۲	۲۱	۹۲/۴۳	۷/۵۷
	Ranker	Info-Gain-Attribute-Eval	۰/۸۲۳	۵۵/۸۱	۱۳	۸۶/۴۹	۱۳/۵۱
		Gain-Ratio-Attribute-Eval	۰/۸۲۸	۵۵/۱۳	۲۳	۸۶/۷۵	۱۳/۲۵
		ReliefF-Attribute-Eval	۰/۸۳۸	۵۲/۸۱	۱۳	۸۷/۵۶	۱۲/۴۴
		Correlation-Attribute-Eval	۰/۸۴۸	۵۲/۶۲	۲۴	۸۸/۳۷	۱۱/۶۳
		Symmetrical-Attribute-Eval	۰/۸۴۳	۵۲/۳۹	۱۹	۸۸/۱۰	۱۱/۹
-	-	۰/۸۳۱	۵۴/۸۴	۴۶	۸۷/۰۲	۱۲/۹۸	
Multilayer Perceptron	Best-First	Cfs-Subset-Eval	۰/۸۳۴	۴۷/۹۳	۱۷	۸۷/۲۹	۱۲/۷۱
	Greedy-Stepwise	Cfs-Subset-Eval	۰/۸۶۵	۴۴/۸۱	۲۱	۸۹/۷۲	۱۰/۲۸
	Ranker	Info-Gain-Attribute-Eval	۰/۷۴۲	۶۱/۹۶	۶	۸۰/۲۷	۱۹/۷۳
		Gain-Ratio-Attribute-Eval	۰/۸۱۷	۵۵/۰۳	۸	۸۵/۹۴	۱۴/۰۶
		ReliefF-Attribute-Eval	۰/۸۷۳	۴۴/۲۰	۴۱	۹۰/۲۷	۹/۷۳
		Correlation-Attribute-Eval	۰/۷۷۰	۵۷/۷۴	۵	۸۲/۴۳	۱۷/۵۷
		Symmetrical-Attribute-Eval	۰/۷۹۵	۵۶/۸۱	۴	۸۴/۳۲	۱۵/۶۸
-	-	۰/۸۸۱	۴۲/۰۹	۴۶	۹۰/۸۱	۹/۱۹	
K-nearest neighbors	Best-First	Cfs-Subset-Eval	۰/۸۸۲	۴۴/۳۰	۱۷	۹۱/۰۸	۸/۹۲
	Greedy-Stepwise	Cfs-Subset-Eval	۰/۸۷۹	۴۴/۰۱	۲۱	۹۰/۸۱	۹/۱۹
	Ranker	Info-Gain-Attribute-Eval	۰/۷۸۵	۵۶/۳۳	۱۵	۸۳/۵۱	۱۶/۴۹
		Gain-Ratio-Attribute-Eval	۰/۸۸۱	۴۸/۲۸	۱۳	۹۱/۰۸	۸/۹۲
		ReliefF-Attribute-Eval	۰/۸۷۱	۴۷/۹۳	۱۳	۹۰/۲۷	۹/۷۳
		Correlation-Attribute-Eval	۰/۸۷۵	۴۵/۲۶	۲۴	۹۰/۵۴	۹/۴۶
		Symmetrical-Attribute-Eval	۰/۸۷۱	۴۶/۸۵	۲۶	۹۰/۲۷	۹/۷۳
-	-	۰/۸۷۷	۴۴/۶۸	۴۶	۹۰/۵۴	۹/۴۶	
Random Forest	Best-First	Cfs-Subset-Eval	۰/۹۰۰	۴۳/۶۵	۱۷	۹۲/۴۳	۷/۷۵
	Greedy-Stepwise	Cfs-Subset-Eval	۰/۹۰۳	۴۴/۳۶	۲۱	۹۲/۷۰	۷/۳۰
	Ranker	Info-Gain-Attribute-Eval	۰/۹۲۱	۴۴/۲۷	۲۸	۹۴/۰۵	۵/۹۵
		Gain-Ratio-Attribute-Eval	۰/۹۱۴	۴۴/۸۶	۴۰	۹۳/۵۱	۶/۴۹
		ReliefF-Attribute-Eval	۰/۹۱۸	۴۴/۷۰	۳۰	۹۳/۷۸	۶/۲۲
		Correlation-Attribute-Eval	۰/۸۹۷	۴۷/۰۴	۲۴	۹۲/۱۶	۷/۸۴
		Symmetrical-Attribute-Eval	۰/۹۱۱	۴۵/۱۲	۳۹	۹۳/۲۴	۶/۷۶
-	-	۰/۹۰۸	۴۷/۳۲	۴۶	۹۲/۹۷	۷/۰۳	

ویژگی‌های بهینه انتخاب شده توسط برخی از الگوریتم‌های مذکور در (جدول ۵) بیان شده است. از این ویژگی‌ها می‌توان به جای تمامی ۶۶ ویژگی برای شناسایی گونه‌ها با کارایی بیشتر استفاده کرد، زیرا به دلیل تعداد ویژگی کمتر و دقت بیشتر در دسته‌بندی (بجز دسته‌بندی MLP) با هزینه بسیار کمتر می‌توان گونه‌ها را شناسایی کرد. سه ویژگی رنگ گلبرگ، تعداد برگچه و نوع کرک در بین الگوریتم‌ها مشترک هستند و ویژگی شکل خار شاخه گل دهنده توسط اکثر الگوریتم‌ها انتخاب شده است. بنابراین چهار ویژگی فوق از پتانسیل بالایی برای انتخاب شدن به عنوان کلید شناسایی بین گونه‌های مختلف تمشک برخوردارند. تمامی این ویژگی‌ها کیفی هستند و جهت بدست آوردن آنها نیازی به هزینه اندازه‌گیری آزمایشگاهی نمی‌باشد.

در مجموع با توجه به (جدول ۴) بیشترین دقت در دسته‌بندی با الگوریتم Random Forest حاصل شد، بطوری که الگوریتم فوق به دقت ۹۴/۰۵ درصد در دسته‌بندی با استفاده از ۲۸ ویژگی انتخاب شده توسط روش جستجوی Ranker و مدل ارزیابی Info Gain Attribute Eval رسید. همچنین با تمامی ۶۶ ویژگی و بدون استفاده از هیچ یک از الگوریتم‌های انتخاب ویژگی، به دقت ۹۲/۹۷ درصد در دسته‌بندی رسید. نکته قابل توجه دیگر این است که دسته‌بندی با الگوریتم MLP با استفاده از چهار ویژگی انتخاب شده توسط روش جستجوی Ranker و مدل ارزیابی Symmetrical Attribute Eval به دقت ۸۴/۳۲ درصد و دسته‌بندی با الگوریتم C4.5 با استفاده از شش ویژگی انتخاب شده توسط روش جستجوی Ranker و مدل ارزیابی Info Gain Attribute Eval به دقت ۸۵/۹۴ درصد رسید.

جدول ۵- ویژگی‌های بهینه انتخاب شده

الگوریتم دسته‌بندی	الگوریتم انتخاب ویژگی		دقت دسته‌بندی (%)	تعداد ویژگی انتخاب شده	ویژگی‌ها	
	روش جستجو	مدل ارزیابی			کیفی	کمی
Random Forest	Ranker	Info Gain Attribute Eval	۹۴/۰۵	۲۸	شکل خار شاخه گل دهنده، رنگ گلبرگ، تعداد برگچه، نوع کرک، فرم ساقه، نوک کاسبرگ، رنگ سطح تحتانی برگچه، کرک بساک، پوشش کرک در سطح فوقانی برگچه، رنگ میله، شکل گوشوارک، نوع کرک کاسبرگ، نوع کرک دمبرگ، نوع گل آذین، نوع کرک سطح فوقانی، نوک برگچه انتهایی، فرم رویشی، حالت شاخه گل دهنده، پوشش کرک در سطح تحتانی، غده کاسبرگ، شکل خار دمبرگ، نوک برگچه جانبی، غده در گل آذین، حضور غده در ساقه	طول دمبرگ، طول برگچه جانبی، طول دمبرگچه انتهایی، طول برگچه انتهایی
K-nearest neighbors	Ranker	Gain Ratio Attribute Eval	۹۱/۰۸	۱۳	نوع گل آذین، فرم رویشی، حالت شاخه گل دهنده، رنگ گلبرگ، نوع کرک کاسبرگ، حضور غده در ساقه، نوع کرک، نوک کاسبرگ، کرک بساک، تعداد برگچه، رنگ میله، غده کاسبرگ، غده در گل آذین	-
C4.5	Ranker	Info Gain Attribute Eval	۸۵/۹۴	۶	رنگ گلبرگ، تعداد برگچه، نوع کرک، شکل خار	طول دمبرگ

	شاخه گل دهنده، فرم ساقه					
-	رنگ گلبرگ، تعداد برگچه، نوع کرک، شکل خار شاخه گل دهنده	۴	۸۴/۳۲	Symmetrical Attribute Eval	Ranker	MLP

References

- [1] Aalders, L. E. and Hall I. V. 1966. A Cytotaxonomic survey of the native blackberries of Nova Scotia. *Canadian Journal of Genetics and Cytology* 8: 528-532.
- [2] Ali, A.S.O., Malik, A.S. and Aziz, A. 2013. A geometrical approach for age-invariant face recognition. *International Visual Informatics Conference*. - Springer 81-96.
- [3] Ballington, JR. Luteyn, MM. Thompson, K. Romoleroux, K. and Castillo, R. 1993. *Rubus* and Vacciniaceous germplasm resources in the Andes of Ecuador. *Plant Genetic Resources newsletter* 93: 9-15.
- [4] Bramer, M. 2007. *Principle of data mining*. Springer.
- [5] Carletta, Jean. 1996. Assessing agreement on classification tasks: The kappa statistic. *Computational Linguistics* 22: 249-254.
- [6] Cohen, J. 1960. A Coefficient of Agreement for Nominal Scales. *Educational and psychological measurement*, 20: 37-46.
- [7] Dash, M. and Liu, H. 1997. Feature selection for classification. *Intelligent data analysis* 1:131-156.
- [8] Diaz, F., & Jones, R. 2004. Using temporal profiles of queries for precision prediction. *International ACM SIGIR conference on Research and development in information retrieval* 18-24.
- [9] Focke, W. O. 1910. *Species Ruborum. Monographiae Generic. Rubi Prodrromus. Bibliotheca Botanica* 17: 1-120.
- [10] Focke W. O. 1911. *Species Ruborum. Monographic Genesis. Rubi Prodrromus. Pars I, Stuttgart*.
- [11] Focke, W. O. 1914. *Species Ruborum. Monographic Genesis. Rubi Prodrromus. Pars I-II. Stuttgart*.
- [12] Gu, Y., C. M. Zhao, W. Jin, and W. L. Li. 1993. *Rubus* resources in Fujian and Hunan provinces. *Acta Horticulturae* 345: 117-125.
- [13] Gustafsson, A. 1942. The origin and properties of the European blackberry flora. *Hereditas* 28: 249-277.
- [14] Gardner, M. W., Dorling, S. R. 1998. Artificial neural networks (the multilayer perceptron)—a review of applications in the atmospheric sciences. *Atmospheric environment*, 32(14-15), 2627-2636.
- [15] Han, J., Pei, J. and Kamber, M. 2011. *Data mining: concepts and techniques*. Elsevier.
- [16] Hummer, K. E. 1996. *Rubus* diversity. *Hort Science* 31: 182-183.
- [17] Ian H. Witten and Eibe Frank. 2005. *Data Mining Practical Machine Learning Tools and Techniques*.
- [18] Kasalkhe, R., Jorjani, E., Sabori, H., Sattarian, A., Habibi, M. 2016. *Biosystematic study of Rubus L. (Rosaceae) in North of Iran*. MSc thesis. University of Gonbad-e-Kavous, 276 pp.
- [19] Kantardzic, M. 2003. *Data Mining: Concepts, models, methods, and algorithms*. Wiley-Interscience.
- [20] Khatamsaz, M. 1992. *Flora of Iran (Rosacea)*.-Research Institute of forests and Rangelands 6: 20-35.
- [21] Larose, D. T. 2005. k-nearest neighbor algorithm. *Discovering Knowledge in Data: An Introduction to Data Mining*, 90-106.
- [22] Lu, L. T. 1983. A study on the genus *Rubus* of china. *Actaphyto taxonomic sinica* 21: 13-25.
- [23] Published on the Internet <https://weka.waikato.ac.nz/explorer> [accessed 13 November 2017].
- [24] Published on the Internet <http://weka.sourceforge.net/doc.dev/weka/attributeSelection/AttributeEvaluator.html> [accessed 13 November 2017].
- [25] Platt, J. 1998. *Sequential Minimal Optimization: A Fast Algorithm for Training Support Vector Machines*. Microsoft Research.
- [26] Quinlan, J.R. 2014. *C4.5: programs for machine learning*. Elsevier 58-60.
- [27] Richards, A. J., J. Kirschner, J. Stepanek, and K. Marhold. 1996. Apomixis and taxonomy: an introduction. *Folia Geobotanica phytotaxonomica* 31: 281-282.

- [28] Rish, I. 2001. An empirical study of the naive Bayes classifier. IJCAI Workshop.
- [29] Robertson, K. R. 1974. The genera of Rosaceae in the southern United States. *Journal of the Arnold Arboretum* 55: 352-360.
- [30] Svetnik, V., Liaw, A., Tong, C., Culberson, J. C., Sheridan, R. P., Feuston, B. P. 2003. Random forest: a classification and regression tool for compound classification and QSAR modeling. *Journal of chemical information and computer sciences*, 43(6), 1947-1958.
- [31] Thompson, M. M. 1995. Chromosome number of *Rubus* species at the National Clonal Germplasm Repository. *Hort Science* 30: 1447-1452.
- [32] Weber, H. E. 1995. Die Gattung *Rubus* L. im nordwestlichen Europa. *Phanerogamarum Monographiae Tomus VII*. J. Cramer, Lehre, Germany.
- [33] Remagnino, P., Mayo, S., Wilkin, P., Cope, J. and Kirkup, D., 2016. *Computational Botany*. Springer Berlin Heidelberg:.
- [34] Ramírez-Gallego, S., García, S., Mouriño-Talín, H., Martínez-Rego, D., Bolón-Canedo, V., Alonso-Betanzos, A., Benítez, J.M. and Herrera, F., 2016. Data discretization: taxonomy and big data challenge. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 6(1), pp.5-21.