

Classification of Web Pages Based on Search Engine Optimization Quality using Random Forest Algorithm

¹Mohammad Hosein Shariatipour*, ²Mohammad Rahmanimanesh, ³Mohammad Javad Parseh

¹Masters student, Computer Engineering, Semnan University, Semnan
mhshariatipour@semnan.ac.ir

²Assistant Professor, Computer Engineering, Semnan University, Semnan
rahmanimanesh@semnan.ac.ir

³Instructor, Computer Engineering, Jahrom University, Jahrom
parseh@jahromu.ac.ir

Abstract

In the evening where ranking on the search engine result pages is directly related to increasing the number of visitors and the progress and development of a business, search engine optimization or SEO is a process that helps to gain a higher ranking. Websites can be classified with the help of machine learning techniques based on the quality of setting SEO guidelines. Classification algorithms are combined with each other with the aim of increasing classification accuracy and are used as an ensemble classification model. In this article, we implement an ensemble classification model with the help of a random forest algorithm, which places web pages in one of the predefined classes based on SEO quality. The obtained results show that the accuracy of the constructed model is between 70.50% and 73.17% and is more accurate than previous works in which ensemble classification algorithms were not used. The built model can help developers build automatic software for detecting the SEO quality of web pages.

Keywords: Search Engine Optimization, Machine Learning, Ensemble Learning, Web Page Classification, Random Forest Algorithm.

طبقه‌بندی صفحه‌های وب بر اساس کیفیت بهینه‌سازی موتور جستجو به کمک الگوریتم جنگل تصادفی

محمدحسین شریعتی‌پور^{۱*}، محمد رحمانی‌منش^۲، محمدجواد پارسه^۳

^۱ دانشجوی کارشناسی ارشد، مهندسی کامپیوتر، دانشگاه سمنان، سمنان

mhshariatipour@semnan.ac.ir

^۲ استادیار، مهندسی کامپیوتر، دانشگاه سمنان، سمنان

rahmanimanesh@semnan.ac.ir

^۳ مربی، مهندسی کامپیوتر، دانشگاه چهرم، چهرم

parseh@jahromu.ac.ir

چکیده

در عصری که رتبه‌بندی در صفحات نتایج موتور جستجو ارتباط مستقیمی با افزایش تعداد بازدیدکنندگان و در نتیجه آن پیشرفت و توسعه یک کسب‌وکار دارد، بهینه‌سازی موتور جستجو یا سئو فرآیندی است که به کسب رتبه بالاتر کمک می‌کند. وبسایت‌ها را می‌توان به کمک تکنیک‌های یادگیری ماشین بر اساس کیفیت تنظیم دستورالعمل‌های سئو طبقه‌بندی کرد. الگوریتم‌های طبقه‌بندی باهدف افزایش دقت طبقه‌بندی با یکدیگر ترکیب می‌شوند و به‌عنوان یک مدل طبقه‌بندی ترکیبی استفاده می‌شوند. در این مقاله یک مدل طبقه‌بندی ترکیبی را به کمک الگوریتم جنگل تصادفی پیاده‌سازی می‌کنیم که صفحات وب را در یکی از طبقه‌بندی‌های از پیش تعریف‌شده بر اساس کیفیت سئو قرار می‌دهد. نتایج به‌دست‌آمده نشان می‌دهد که دقت مدل ساخته‌شده بین ۷۰/۵۰٪ تا ۷۳/۱۳٪ است و نسبت به کارهای قبلی که در آن‌ها از الگوریتم‌های طبقه‌بندی ترکیبی استفاده‌نشده است دقت بالاتری دارد. مدل ساخته‌شده می‌تواند به توسعه‌دهندگان نرم‌افزار برای ساخت نرم‌افزارهای خودکار تشخیص کیفیت سئو صفحات وب کمک کند.

کلمات کلیدی

بهینه‌سازی موتور جستجو، یادگیری ماشین، یادگیری ترکیبی، طبقه‌بندی صفحات وب، الگوریتم جنگل تصادفی

جستجو مربوط می‌شود [4]. در عصر حاضر، بهینه‌سازی موتور جستجو برای

کسب‌وکارها امری ضروری است تا بتوانند در فضای رقابتی وب باقی بمانند.

مطالعات بسیاری در ارتباط با فاکتورهای بهینه‌سازی موتور جستجو و تأثیر آن بر کسب رتبه بالاتر در نتایج موتور جستجو انجام گرفته است [1, 4-6]. به‌طور کلی این فاکتورها به دو دسته درون صفحه و خارج صفحه تقسیم‌بندی می‌شوند. عنوان صفحه، بهینه‌سازی فایل‌های گرافیکی، نام دامنه، استفاده از متاتگ‌ها و کلیه فاکتورهایی که تحت کنترل کامل مدیر وبسایت قرار دارند، فاکتورهای درون صفحه و ساخت یک لینک‌ها، فعالیت در رسانه‌های اجتماعی و به‌طور کلی عواملی که مدیر وبسایت کنترل کمتری بر روی آن‌ها دارد فاکتورهای خارج از صفحه در فرآیند بهینه‌سازی موتور جستجو محسوب می‌شوند.

امروزه طبقه‌بندی^۲ صفحات وب به‌منظور اهداف مختلفی به کار می‌رود، طبقه‌بندی صفحات وب بر اساس کیفیت سئو یکی از موضوعاتی است که اخیراً مورد توجه قرار گرفته است. الگوریتم‌های طبقه‌بندی ویژگی‌های متفاوتی

۱- مقدمه

درحالی‌که دیجیتالی شدن کسب‌وکارها در حال رشد و پیشرفت است، از وبسایت‌ها می‌توان به‌عنوان ابزاری برای گسترش آگاهی و توجه به خدمات خود در شبکه اینترنت استفاده کرد [1]. بیش از نیمی از ترافیک یک وبسایت توسط موتورهای جستجو به آن هدایت می‌شود و سهم گوگل از بازار موتورهای جستجو تقریباً ۹۲/۴۲٪ است، به همین دلیل رتبه‌بندی در موتور جستجو گوگل از اهمیت ویژه‌ای برخوردار است [2].

اغلب کاربران فقط از صفحاتی که در صفحه اول نتایج موتور جستجو ارائه می‌شود بازدید می‌کنند و نتایج در صفحات بعدی شانس کمتری برای بازدید دارند بنابراین، صفحات وب با رتبه بالا بازدیدکنندگان بیشتری داشته و همین امر، بالا بردن رتبه‌بندی صفحات را به اولویت اصلی مدیران وب تبدیل کرده است [3]. بهینه‌سازی موتور جستجو^۱ یا سئو به فعالیت بهینه‌سازی وبسایت‌ها و صفحات وب برای کسب رتبه بالاتر در صفحات نتایج موتور

عملکرد الگوریتم‌های طبقه‌بندی یادگیری ماشینی؛ k نزدیک‌ترین همسایه؛ ماشین بردار پشتیبان و جنگل تصادفی برای تشخیص صفحات فیشینگ در پژوهش [10] ارزیابی شده است، این الگوریتم‌ها بر روی مجموعه داده‌های مشکل از ۱۰۰۰ صفحه فیشینگ و ۴۰۰ صفحه قانونی که هر یک از صفحات شامل ۱۴ ویژگی است پیاده‌سازی شد و نتایج نشان داد الگوریتم طبقه‌بندی جنگل تصادفی با دقت ۹۸/۳۵٪ نسبت به دو مدل دیگر عملکرد بهتری دارد.

از رویکردهای یادگیری ماشینی برای شناسایی کیفیت سئو یک صفحه وب در [11] استفاده شد و ۵ الگوریتم طبقه‌بندی درخت تصمیم، k نزدیک‌ترین همسایه، ماشین بردار پشتیبان، طبقه‌بندی ساده بیز و رگرسیون لجستیک^۲ بر روی مجموعه‌ای از ۶۰۰ صفحه وب که کیفیت سئو آن‌ها توسط متخصصان سئو مشخص شده بود آموزش داده شدند. ارزیابی‌ها نشان داد این الگوریتم‌های طبقه‌بندی با دقت ۵۴/۶۹٪ تا ۶۹/۶۷٪ نسبت به دقت پایه که مربوط به طبقه‌بند اکثریت (۴۸/۸۳٪) است، دقت بالاتری دارند.

یک مدل طبقه‌بندی صفحات وب بر اساس ترکیبی از طبقه‌بندی کننده‌های متعدد در مقاله [12] پیشنهاد شده است که در آن از ساختار درخت مانند تگ‌های HTML برای مشخص کردن ویژگی‌های ساختاری یک صفحه وب استفاده می‌شود، این ویژگی‌های ساختاری به همراه ویژگی‌های متنی صفحه به بردار تبدیل شده و باهم ترکیب می‌شوند. ارزیابی مدل پیشنهادی بر روی مجموعه داده‌های عمومی موجود به دقت قابل قبولی رسیده است و نتایج اثبات می‌کند دقت طبقه‌بندی صفحات وب با ترکیب چند طبقه‌بند متعدد بهبود می‌یابد.

۳- یادگیری ترکیبی

یادگیری ترکیبی یک الگوی موفق از یادگیری ماشینی است که مجموعه‌ای از یادگیرندگان را به جای استفاده از یک یادگیرنده واحد برای پیش‌بینی ویژگی نمونه جدید ادغام می‌کند. در این ساختار، تمام مقادیر خروجی به دست آمده از هر یادگیرنده با استفاده از یک مکانیسم رأی‌گیری برای پیش‌بینی برچسب کلاس نهایی ترکیب می‌شود. هدف اصلی یادگیری ترکیبی تشکیل یک طبقه‌بندی کننده قوی مشکل از چند یادگیرنده برای به دست آوردن نتایج طبقه‌بندی دقیق‌تر است [13].

۱- ۳- Bagging

Bagging یکی از پرکاربردترین روش‌های یادگیری ترکیبی است که از روش نمونه‌گیری مجدد Bootstrapping استفاده می‌کند. در این روش مجموعه آموزشی اولیه با چند زیرمجموعه از داده‌های اولیه (آزمایش‌های Bootstrapping) در طول فرآیند Bagging جایگزین می‌شود. در هر یک از این زیرمجموعه‌ها، برخی از مقادیر داده‌ها می‌توانند مکرر ظاهر شوند این در حالی است که برخی دیگر ممکن است در هیچ‌یک از زیرمجموعه‌ها عضو نباشند. سپس نتیجه نهایی با میانگین‌گیری نتایج هر زیرمجموعه تعیین می‌شود [14].

دارند و عملکرد آن‌ها تحت تأثیر عوامل مختلفی از جمله مجموعه داده‌ها است. همچنین الگوریتم‌های طبقه‌بندی را می‌توان با یکدیگر ترکیب کرده و به عنوان یک مدل طبقه‌بندی ترکیبی از آن استفاده کرد.

در پژوهش پیش رو از فاکتورهای درون صفحه و تکنیک‌های یادگیری ماشینی برای ساخت یک مدل طبقه‌بندی استفاده می‌شود که به طور خودکار صفحات وب را بر اساس میزان تنظیم دستورالعمل‌های سئو طبقه‌بندی کند. به طور ویژه هدف در این پژوهش استفاده از الگوریتم طبقه‌بندی جنگل تصادفی^۳ به عنوان یک الگوریتم یادگیری ترکیبی و مقایسه دقت به دست آمده با دقت سایر الگوریتم‌های طبقه‌بندی که در پژوهش‌های پیشین استفاده شده است خواهد بود. از مدل طبقه‌بندی ترکیبی ایجاد شده می‌توان به عنوان ابزاری برای تشخیص صفحاتی که نیاز به بهینه‌سازی موتور جستجو دارند استفاده کرد.

ادامه این مقاله به شرح زیر سازمان دهی شده است. در بخش ۲ پیشنهادی پژوهش را بررسی خواهیم کرد. بخش ۳ مفهوم یادگیری ترکیبی^۴ را تشریح می‌کند. بخش ۴ به معرفی مجموعه داده‌ها می‌پردازد. بخش ۵ روش پیشنهادی را توضیح می‌دهد که چگونه یک مدل طبقه‌بندی ترکیبی مبتنی بر الگوریتم جنگل تصادفی پیاده‌سازی کنیم تا به کمک آن صفحات وب را بر اساس کیفیت سئو طبقه‌بندی کند. نتایج ارزیابی نیز در بخش ۶ ارائه می‌شود.

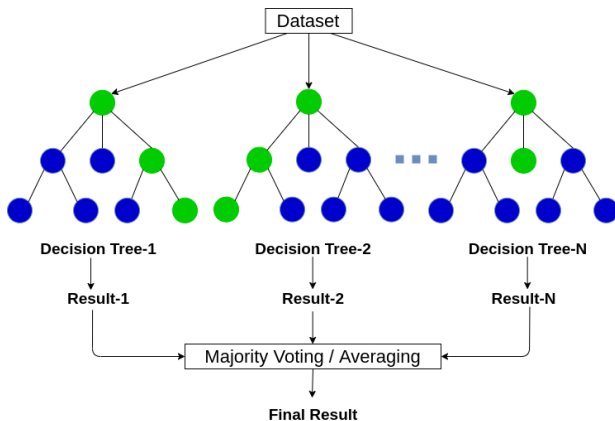
۲- کارهای مرتبط

در [7] نویسندگان عوامل تأثیرگذار سئو بر رتبه‌بندی نتایج موتور جستجو را در پژوهش‌های پیشین مطالعه کرده و ۲۴ عامل را به عنوان مهم‌ترین عوامل معرفی کرده‌اند. از جمله این موارد کیفیت بک لینک، کلمه کلیدی در تگ عنوان سایت، زمان بارگذاری سایت و چگالی کلمات کلیدی است.

استفاده از روش‌های تصمیم‌گیری برای بررسی رابطه متقابل و وزن‌های تأثیرگذار در میان معیارهای سئو در [8] نشان می‌دهد متاگ‌ها مهم‌ترین معیار در سئو هستند و پس از آن کلمات کلیدی، طراحی وب‌سایت، لینک سازی و رسانه‌های اجتماعی قرار دارند. کارشناسان پیشنهاد می‌کنند که مدیران وب‌سایت‌ها بیشترین تأکید را بر متاگ‌ها داشته باشند.

طبقه‌بندی صفحات وب برای اهداف مختلف یک زمینه مطالعاتی گسترده در چند سال اخیر بوده است.

وب‌سایت‌های بدافزار و کرک اغلب آدرس اینترنتی خود را برای جلوگیری از شناسایی خودکار تغییر می‌دهند. با این حال، در بسیاری از موارد، این وب‌سایت‌ها طرح‌های بصری خاصی مانند رنگ‌های خاص، فونت‌ها، اشکال و اندازه‌ها را حفظ می‌کنند که می‌تواند معرف عملکرد وب‌سایت به کاربران باشد. نویسندگان در مقاله [9] از ویژگی‌های طراحی بصری و غیر بصری وب‌سایت‌ها برای شناسایی دسته‌بندی وب‌سایت، به ویژه وب‌سایت‌های بدافزار و کرک استفاده کردند و وب‌سایت‌های مخرب با نرخ مثبت واقعی نسبتاً بالا و نرخ مثبت کاذب ناچیزی شناسایی شدند. نتایج نشان می‌دهد استفاده از روش‌های طبقه‌بندی که بر ویژگی‌های طراحی متکی است، می‌تواند به عنوان یک جایگزین مؤثر برای متن کاوی در دسته‌بندی وب‌سایت‌ها باشد که اغلب فرآیندی زمان‌بر و نیاز به محاسبات پیچیده دارد.



شکل (۱): الگوریتم جنگل تصادفی

جدول (۱): فاکتورهای مؤثر بر SEO داخلی

عنوان	توضیحات
Tlen	طول محتوا در تگ <title> (تعداد کلمات)
Tkw	فراوانی کلمات کلیدی در تگ <title>
Mlen	طول محتوا در تگ <meta> توضیحات
Mkw	فراوانی کلمات کلیدی در تگ <meta> توضیحات
h1	تعداد تگ <h1>
h1len	میانگین طول محتوا تگ‌های <h1>
h1kw	فراوانی کلمات کلیدی در تگ‌های <h1>
h2	تعداد تگ <h2>
h2len	میانگین طول محتوا تگ‌های <h2>
h2kw	فراوانی کلمات کلیدی در تگ‌های <h2>
h3	تعداد تگ <h2>
h3len	میانگین طول محتوا تگ‌های <h2>
h3kw	فراوانی کلمات کلیدی در تگ‌های <h2>
h3kw	تعداد تگ‌های که حاوی ویژگی alt هستند
altkw	فراوانی کلمات کلیدی در ویژگی alt تگ‌های
linkkw	تعداد فراوانی کلمه کلیدی در متن anchor
linkout	تعداد پیوندهای خروجی
urlen	طول آدرس اینترنتی (url)
urlkw	فراوانی کلمات کلیدی در آدرس اینترنتی (url)
txtlen	طول متن بدنه صفحه
txtkw	تعداد کلمات کلیدی در متن بدنه صفحه

۱-۱-۳- جنگل تصادفی

الگوریتم جنگل تصادفی، نمونه‌ای از یک الگوریتم Bagging است که جنگلی با درخت‌های تصمیم‌گیری متعدد ایجاد می‌کند. ایده اصلی الگوریتم جنگل تصادفی استفاده از روش نمونه‌برداری Bootstrap برای انجام عملیات نمونه‌برداری با جایگزینی از مجموعه داده اصلی و سپس ساختن درخت تصمیم از زیرمجموعه داده‌های نمونه‌برداری شده و ترکیب درخت‌های تصمیم متعدد در یک جنگل تصادفی است [15].

هر درخت تصمیم در الگوریتم جنگل تصادفی مطابق شکل (۱) شامل یک دنباله درخت از گره‌های تصمیم است. بر اساس این توالی، درخت به شاخه‌های مختلف تقسیم می‌شود تا به برگ‌های درخت برسد. نتایج پیش‌بینی هر درخت تصمیم از طریق گره‌های برگ مشخص می‌شود و در نهایت خروجی‌های درخت‌های تصمیم چندگانه برای پیش‌بینی ترکیب می‌شوند و بر اساس اکثریت نتایج به‌دست‌آمده از درخت‌های تصمیم چندگانه، پیش‌بینی نهایی جنگل تصادفی نتیجه می‌شود.

وجود درختان تصمیم متفاوت در الگوریتم جنگل تصادفی به کاهش مشکل بیش‌برازش کمک می‌کند، مشکلی که مخصوص الگوریتم درخت تصمیم است [10].

۴- مجموعه داده‌ها

در کار پیش رو از مجموعه داده‌های [16] جمع‌آوری شده در پژوهش [11] استفاده خواهیم کرد. این مجموعه داده شامل ۶۰۰ نمونه تصادفی از صفحات وب است که هر یک از صفحات توسط ۲۱ متغیر مستقل به‌عنوان فاکتورهای مؤثر بر سئو داخلی که در جدول (۱) قابل مشاهده است و یک متغیر وابسته به‌عنوان کیفیت سئو صفحه توصیف می‌شود. کیفیت سئو صفحات توسط متخصصان بهینه‌سازی موتور جستجو در یکی از دسته‌های ضعیف، متوسط و قوی قرار گرفته است.

۵- روش پژوهش

هدف از این پژوهش ساخت مدل طبقه‌بندی کننده ترکیبی است که صفحات وب را بر اساس کیفیت سئو در یکی از دسته‌های قوی، متوسط و ضعیف قرار دهد. از الگوریتم جنگل تصادفی به‌عنوان یک الگوریتم یادگیری ترکیبی از نوع Bagging برای ساخت مدل پیشنهادی استفاده خواهیم کرد. مراحل ساخت مدل طبقه‌بندی و ارزیابی آن در ۴ مرحله دنبال می‌شود.

۱-۵- نرمال‌سازی داده‌ها

نرمال‌سازی معمولاً فرآیندی است که باهدف بهبود دقت و عملکرد مدل‌های طبقه‌بندی بر روی داده‌های عددی انجام می‌شود تا مجموعه داده به یک بازه خاص، معمولاً بین ۰ و ۱ محدود شود. نرمال‌سازی در مواردی که داده‌های عددی دارای محدوده‌های مختلف باشند ضروری است. از آنجایی که کلیه ۲۱ ویژگی مستقل موجود در مجموعه داده‌ها از نوع عددی هستند ابتدا لازم است مقادیر را نرمال کنیم.

```

from sklearn.ensemble import RandomForestClassifier

clf_model=RandomForestClassifier(
    n_estimators = 100,
    min_samples_leaf = 16
)
    
```

شکل (۲): تعریف مدل طبقه‌بندی جنگل تصادفی

اعتبارسنجی Hold-out تکنیکی است که در آن مجموعه داده‌ها به دو بخش آموزش و آزمون تقسیم می‌شود. مجموعه آموزش برای آموزش مدل یادگیری و مجموعه آزمون برای اعتبارسنجی عملکرد مدل ایجاد شده. اما این تکنیک در مواردی که مجموعه داده کوچک باشد و یا اگر مجموعه داده آموزش انتخاب شده معرف خوبی برای کل داده‌ها نباشد ممکن است عملکرد مطلوبی نداشته باشد.

اعتبارسنجی متقاطع K-fold تکنیک دیگری است که در یادگیری ماشین برای ارزیابی عملکرد یک مدل استفاده می‌شود و محدودیت‌های تکنیک قبل را برطرف می‌کند. در این روش مجموعه داده‌ها به k قسمت مساوی تقسیم می‌شود. سپس مدل یادگیری بر روی K-1 قسمت مرحله آموزش را طی می‌کند و پس از آن با یک قسمت باقیمانده اعتبارسنجی مدل سنجیده می‌شود. این فرآیند k بار تکرار می‌شود و هر بار از یک قسمت متفاوت به عنوان مجموعه آزمون استفاده می‌کنیم. در آخر با میانگین‌گیری از k مرحله انجام شده ارزیابی انجام خواهد شد.

از معیار دقت نیز برای ارزیابی مدل طبقه‌بندی استفاده خواهیم کرد، این معیار نشان می‌دهد چه تعداد از نمونه‌ها در مرحله آزمون به درستی طبقه‌بندی شده‌اند. برای محاسبه دقت، تعداد داده‌های درست طبقه‌بندی شده بر تعداد کل داده‌ها تقسیم می‌شود.

۶- نتایج

ارزیابی مدل طبقه‌بندی توسط هر دو تکنیک Hold-out و K-fold انجام خواهد شد. در تکنیک Hold-out از دوسوم داده‌ها برای مرحله آموزش و از یک‌سوم باقیمانده برای آزمون استفاده می‌شود. در روش K-fold نیز مقدار k را برابر ۱۰ در نظر می‌گیریم.

دقت ۵ الگوریتم طبقه‌بندی درخت تصمیم، k نزدیک‌ترین همسایه، ماشین بردار پشتیبان، طبقه‌بندی ساده بیز و رگرسیون لجستیک که در پژوهش [۱۱] به دست آمده است به همراه الگوریتم طبقه‌بندی ترکیبی جنگل تصادفی ارائه شده در این پژوهش به دو روش Hold-out و K-fold به ترتیب در جدول (۲) و جدول (۳) آورده شده.

حداکثر دقت به دست آمده در پژوهش [11] در حالت Hold-out مربوط به الگوریتم طبقه‌بندی درخت تصمیم با مقدار ۶۷/۶۵٪ و در حالت K-fold حداکثر دقت مربوط به الگوریتم k نزدیک‌ترین همسایه با مقدار ۶۷/۶۹٪ درصد است، این در حالی است که دقت‌های به دست آمده در مدل پیشنهادی ما به ترتیب برابر ۷۰/۵۰٪ و ۷۳/۱۷٪ است.

نرمال‌سازی MinMax تکنیکی است که از آن استفاده خواهیم کرد و به کمک آن کلیه مقادیر عددی را محدود به بازه ۰ تا ۱ می‌کنیم. نحوه محاسبه مقدار نرمال متغیر X به کمک فرمول (۱) محاسبه می‌شود.

$$X_{norm} = \frac{X - X_{min}}{X_{max} - X_{min}} \quad (1)$$

که در آن، X_{norm} نشان‌دهنده مقدار نرمال شده، X_{min} و X_{max} به ترتیب حداقل و حداکثر مقادیر داده‌ها هستند و X نشان‌دهنده مقدار اصلی داده است.

۲-۵- تنظیم فرآیندها

فرآیندها مقادیری هستند که قبل از شروع فرآیند یادگیری باهدف بهینه‌سازی مدل [17] و تطبیق یک مدل یادگیری ماشین در مسائل مختلف تنظیم می‌شوند [18]. انتخاب فرآیندهای مناسب برای مدل‌های یادگیری ماشین تأثیر مستقیمی بر عملکرد مدل دارد که اغلب به دانش عمیق از الگوریتم‌های یادگیری ماشین و تکنیک‌های بهینه‌سازی فرآیندها نیاز دارد.

تعداد درختان (n_estimators) و حداقل نمونه‌های موردنیاز در یک گرهی برگ (min_sample_leaf) از جمله فرآیندهای پرکاربرد در الگوریتم جنگل تصادفی هستند که از آن‌ها برای ساخت مدل طبقه‌بندی خود استفاده خواهیم کرد.

۳-۵- ساخت مدل طبقه‌بندی

در کار پیش رو از الگوریتم جنگل تصادفی موجود در کتابخانه Scikit Learn برای ساخت مدل طبقه‌بندی استفاده می‌کنیم، Scikit Learn یک کتابخانه پرکاربرد زبان برنامه‌نویسی پایتون برای اهداف یادگیری ماشین است که ابزارهای کاربردی زیادی به منظور یادگیری ماشین و مدل‌سازی آماری داده‌ها همچون طبقه‌بندی و خوشه‌بندی فراهم می‌کند.

مقدار فرآیندها تعداد درختان را در مدل خود برابر ۱۰۰ و حداقل نمونه‌های موردنیاز در یک گرهی برگ را ۱۶ در نظر می‌گیریم، شکل (۲) نحوه تعریف مدل طبقه‌بندی جنگل تصادفی با این مقادیر به زبان پایتون را نشان می‌دهد. مدل تعریف شده را با ۲۱ فاکتور مؤثر بر سؤ داخلی که در مجموعه داده‌ها وجود دارد آموزش می‌دهیم و انتظار داریم این مدل پس از طی کردن مرحله آموزش بتواند کیفیت سؤی نمونه وب‌سایت‌های جدید را با دقت قابل قبولی پیش‌بینی کند.

۴-۵- ارزیابی مدل طبقه‌بندی

از آنجایی که عوامل مختلفی از جمله دامنه داده‌ها، ویژگی آن‌ها و فرضیات الگوریتم بر روی الگوریتم‌های طبقه‌بندی تأثیرگذار است انتخاب بهترین الگوریتم طبقه‌بندی برای یک کار خاص نیاز به ارزیابی الگوریتم‌های مختلف دارد. دو تکنیک رایج برای ارزیابی الگوریتم‌های طبقه‌بندی وجود دارد.

performance and compliance assessment based on a data-driven search engine optimization methodology," *Information*, vol. 12, no. 7, p. 259, 2021.

- [2] M. Vález and A. Ventura, "Analysis of the SEO visibility of university libraries and how they impact the web visibility of their universities," *The Journal of Academic Librarianship*, vol. 46, no. 4, p. 102171, 2020.
- [3] A.-J. Su, Y. C. Hu, A. Kuzmanovic, and C.-K. Koh, "How to improve your search engine ranking: Myths and reality," *ACM Transactions on the Web (TWEB)*, vol. 8, no. 2, pp. 1-25, 2014.
- [4] M. Khan and A. Mahmood, "A distinctive approach to obtain higher page rank through search engine optimization," *Sāhdhanā*, vol. 43, no. 3, p. 43, 2018.
- [5] S. An and J. J. Jung, "A heuristic approach on metadata recommendation for search engine optimization," *Concurrency and Computation: Practice and Experience*, vol. 33, no. 3, p. e5407, 2021.
- [6] A. Erdmann, R. Arilla, and J. M. Ponzoa, "Search engine optimization: The long-term strategy of keyword choice," *Journal of Business Research*, vol. 144, pp. 650-662, 2022.
- [7] C. Ziakis, M. Vlachopoulou, T. Kyrkoudis, and M. Karagkiozidou, "Important factors for improving Google search rank," *Future internet*, vol. 11, no. 2, p. 32, 2019.
- [8] H.-J. Tsuei, W.-H. Tsai, F.-T. Pan, and G.-H. Tzeng, "Improving search engine optimization (SEO) by using hybrid modified MCDM models," *Artificial Intelligence Review*, vol. 53, pp. 1-16, 2020.
- [9] D. Cohen, O. Naim, E. Toch, and I. Ben-Gal, "Website categorization via design attribute learning," *Computers & Security*, vol. 107, p. 102312, 2021.
- [10] T. O. Ojewumi, G. Ogunleye, B. Oguntunde, O. Folorunsho, S. Fashoto, and N. Ogbu, "Performance evaluation of machine learning tools for detection of phishing attacks on web pages," *Scientific African*, vol. 16, p. e01165, 2022.
- [11] G. Matošević, J. Dobša, and D. Mladenčić, "Using machine learning for web page classification in search engine optimization," *Future Internet*, vol. 13, no. 1, p. 9, 2021.
- [12] L. Deng, X. Du, and J.-z. Shen, "Web page classification based on heterogeneous features and a combination of multiple classifiers," *Frontiers of Information Technology & Electronic Engineering*, vol. 21, no. 7, pp. 995-1004, 2020.
- [13] P. Y. Taser, "Application of bagging and boosting approaches using decision tree-based algorithms in diabetes risk prediction," in *Proceedings*, 2021, vol. 74, no. 1: MDPI, p. 6.
- [14] A. Aldrees, H. H. Awan, M. F. Javed, and A. M. Mohamed, "Prediction of water quality indexes with ensemble learners: Bagging and Boosting," *Process Safety and Environmental Protection*, vol. 168, pp. 344-361, 2022.
- [15] Z. Chen, Q. Wu, S. Han, J. Zhang, P. Yang, and X. Liu, "A study on geological structure prediction based on random forest method," *Artificial Intelligence in Geosciences*, vol. 3, pp. 226-236, 2022.

جدول (۲): دقت در تکنیک Hold-out

دقت	فراپارامتر	الگوریتم
۶۵٪/۶۷	C = 0.49 min_samples_leaf=16	Decision trees
۵۸٪/۷۱	-	Naïve Bayes
۶۵٪/۱۷	K = 45, p = 2	KNN
۶۲٪/۶۸	C = 7.74 × 10 ³ σ = 0.000464	SVM
۶۲٪/۱۹	C = 2	Logistic regression
۷۰٪/۵۰	n_estimators=100 min_samples_leaf=16	Random forest

جدول (۳): دقت در تکنیک K-fold

دقت	فراپارامتر	الگوریتم
۶۷٪/۵۳	C = 0.49 min_samples_leaf=16	Decision trees
۵۴٪/۶۹	-	Naïve Bayes
۶۹٪/۶۷	K = 45, p = 2	KNN
۶۶٪/۱۸	C = 7.74 × 10 ³ σ = 0.000464	SVM
۶۲٪/۹۹	C = 2	Logistic regression
۷۳٪/۱۷	n_estimators=100 min_samples_leaf=16	Random forest

از نتایج به دست آمده می توان نتیجه گرفت استفاده از مدل یادگیری ترکیبی می تواند دقت را در طبقه بندی صفحات وب بر اساس کیفیت بهینه سازی موتور جستجو بهبود دهد. از مدل طبقه بندی ترکیبی ساخته شده می توان برای ساخت نرم افزارهای خودکار پیش بینی کیفیت سئو صفحات وب یا به عنوان ابزاری برای شناسایی صفحاتی که نیاز به بهبود کیفیت سئو دارند استفاده کرد.

مجموعه داده های استفاده شده در کار ما شامل ۲۱ عامل مؤثر بر سئو داخلی صفحات بود این در حالی که است که عوامل خارجی نیز تأثیر قابل توجهی بر کیفیت سئو صفحات دارند بنابراین در کارهای آینده لازم است از ترکیب عوامل داخلی و خارجی برای ارزیابی کیفیت سئو صفحات وب استفاده شود. استفاده از سایر الگوریتم های طبقه بندی ترکیبی نیز می تواند یکی دیگر از کارهای پیش رو باشد.

مراجع

- [1] I. Drivas, D. Kouis, D. Kyriaki-Manessi, and G. Giannakopoulos, "Content management systems

- [16] Dataset: <https://zenodo.org/record/4416123>
- [17] D. Sun, H. Wen, D. Wang, and J. Xu, "A random forest model of landslide susceptibility mapping based on hyperparameter optimization using Bayes algorithm," *Geomorphology*, vol. 362, p. 107201, 2020.
- [18] L. Yang and A. Shami, "On hyperparameter optimization of machine learning algorithms: Theory and practice," *Neurocomputing*, vol. 415, pp. 295-316, 2020.

زیرنویس ها

- ¹ Search engine optimization
- ² Classification
- ³ Random Forest
- ⁴ Ensemble Learning
- ⁵ K-Nearest Neighbors
- ⁶ Naive Bayes
- ⁷ Logistic regression
- ⁸ Overfitting