

تشخیص حس متن با استفاده از روش‌های یادگیری ماشین

دکتر امیر رجبی	پری دریاکش	مونا رجبی فرد	کمیل آقابابایی	مریم یوسفی زاده
دانشگاه آزاد بندرعباس	دانشگاه آزاد بندرعباس	دانشگاه آزاد بندرعباس	دانشگاه آزاد بندرعباس	دانشگاه آزاد بندرعباس
A_rajabi@yahoo.com	Pari.daryakesh@yahoo.com	m.rajabi@tci.ir	babaiekomeil@gmail.com	m-yousefizadeh@irimo.ir

چکیده

در متن کاوی به منظور تشخیص احساس متن برای دسته بندی و بررسی نظرها در سایت‌های اجتماعی به منظور تشخیص قطبیت متن، مطالعاتی انجام شده است. هدف این مقاله، بررسی تأثیر الگوریتم‌های طبقه‌بندی برای آنالیز حس متن می‌باشد. در این مقاله، آنالیز احساس به روش نظارت شده مبتنی بر داده‌های برچسب‌گذاری انجام شده است. به کارگیری الگوریتم‌های متفاوت در میزان تشخیص حس جمله می‌تواند بسیار موثر باشد. از این رو این مقاله، با در نظر گرفتن الگوریتم‌های طبقه‌بندی جنگل‌های تصادفی، نایو بیز، درخت اضافی و SVC خطی بررسی و نوشته شده است.

بررسی میزان تأثیر با استفاده از معیارهای Recall و Precision و F-score سنجیده می‌شود. نتایج به دست آمده از اجرای الگوریتم‌ها بر روی داده‌های آزمون، سنجیده خواهند شد. بررسی نتایج نشان می‌دهد انتخاب الگوریتم Random Forest می‌تواند، میزان دقت را افزایش دهد.

واژه‌های کلیدی: یادگیری ماشین، تشخیص احساسات، طبقه‌بند Linear SVC، طبقه‌بند Naive Bayes، طبقه‌بند Random Forest، طبقه‌بند Extra Trees، Bag of Words، TF-IDF

۱. مقدمه

امروزه آنالیز احساس در داده‌های متنی، زمینه مطالعاتی است که سعی در بیان احساس‌ها، رفتارها، نظرها و تحلیل افراد مختلف نسبت به موجودیت‌ها و ویژگی‌های آن دارد. این موجودیت می‌تواند محصول، سرویس، سازمان، فرد، رخداد و موضوع باشد. هدف از آنالیز یا تحلیل احساسات، پیدا کردن نظرهایی است که احساسی را نشان داده و جهت‌گیری این نظرها را تشخیص می‌دهد. با رشد رسانه‌های اجتماعی مانند نظرسنجی‌ها، فروم‌ها، انجمن‌های گفت‌وگو، وبلاگ‌ها، توییتر و شبکه‌های اجتماعی، اهمیت آنالیز احساسات افزایش یافته است. سیستم‌های آنالیز احساسات کمابیش در همه زمینه‌های تجاری و اجتماعی به کار گرفته می‌شوند؛ زیرا نظرها و عقیده‌ها برای همه فعالیت‌های انسانی مهم بوده و تأثیر شایان توجهی بر رفتار ما دارند.

در این مقاله هدف، یافتن نگرش متن نظر با توجه به دیدگاه و موضوعات مطرح شده در متن خبر است. در این راستا مجموعه داده‌ای متشکل از اخبار و نظرات جمع‌آوری شده و توسط افراد خبره برچسب زده شده و روشی پیشنهادی با استفاده از روش‌های متن کاوی ارائه شده است. این روش بر مبنای تحلیل ساختار خبر، یافتن ارتباط بین متن نظر و خبر می‌باشد. این ارتباط از تجزیه ساختار دستوری متن و استخراج اسامی برای یافتن موضوع متن خبر و سپس یافتن اشتراک در متن نظر و خبر بدست می‌آید. ساختار این مقاله به شکل زیر است. پس از مقدمه در بخش ۲ نگاهی بر مطالعات انجام شده در سال‌های اخیر داریم و در بخش ۳ با تشریح مجموعه داده و در بخش ۵ به مقایسه تأثیر الگوریتم‌های طبقه‌بند می‌پردازیم.

۲. پیشینه پژوهش

آنالیز احساس در حال حاضر به موضوع روز تبدیل شده و تحقیقات زیادی در این زمینه انجام شده است. از جمله تحقیقات بر چگونگی استخراج ویژگی و تغییر ویژگی ها برای افزایش دقت الگوریتم انجام شده است.

در سال ۲۰۱۸ Arxiv JiHoPark با استفاده از مدل های یادگیری ماشین سنتی مانند رگرسیون بردار پشتیبانی (SVR) و رگرسیون منطقی به بررسی تحلیل احساسات متن پرداختند.

Christos Baziotis, Arxiv در سال ۲۰۱۸ با استفاده از LSTM دو طرفه با مکانیزم توجه عمیق به پیش بینی محتوای موثر توییت با RNN های عمیق توجه و یادگیری انتقال پرداختند.

UmangGupta Ankush Chatterjee در سال ۲۰۱۸ استفاده از یک مدل یادگیری عمیق مبتنی بر LSTM, تشخیص احساسات در مکالمات متنی را بررسی کردند و به کشف احساسات فیزیولوژیکی در حوزه ی علوم رایانه ای پرداختند.

Ms. Farha Nausheen در سال ۲۰۱۸ با جمع آوری توییت کاربران خاص با استفاده از کتابخانه twitter API Twython در پایتون تجزیه و تحلیل احساسات برای پیش بینی نتایج انتخابات پرداخته است.

در سال ۲۰۱۸ Ji Ho Park Arxiv با استفاده از یک مدل CNN بنام EVEC برای مدلسازی احساس درون یک کلمه است که قابلیت ترکیب با کارهای طبقه بندی متن راداراست به پیدا کردن نماینده های خوب از احساسات برای طبقه بندی متن پرداخته است.

Hakak Nida در سال ۲۰۱۹ با استفاده از مدل SVM به طبقه بندی احساسات خودکار پرداخته است.

Spencer Cappallo در سال ۲۰۱۸ با استفاده از شبکه عصبی پیشرفته LSTM و مجموعه ای بزرگ از مقادیر Emoji برای بیان روابط معنایی بین Emoji و متن/Emoji و تصاویر به چالش های موجود در پیش بینی ایموجی پرداخته است.

Toshiki Tomihira در سال ۲۰۱۸ با استفاده از مقایسه مدل رمزگذار-رمزگشای شبکه عصبی (RNN) و شبکه عصبی (CNN) به پیش بینی ایموجی عصبی برای تحلیل احساسات پرداخته است.

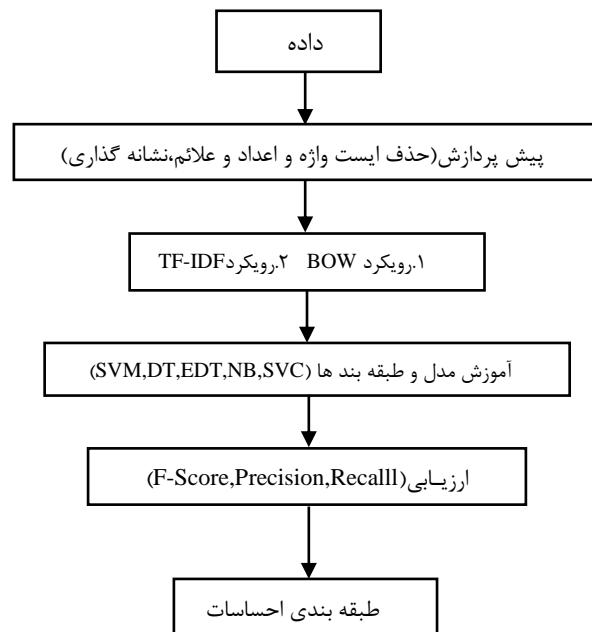
Bharat Gaind در سال ۲۰۱۹ با استفاده از پردازش زبان طبیعی و الگوریتم طبقه بندی Machine Learning به تشخیص احساس و تجزیه و تحلیل در رسانه های اجتماعی پرداخته است.

Kazuyuki Matsumoto Tokushim در سال ۲۰۱۸ با استفاده از دقت بالای روش های مبتنی بر BiLSTM نسبت به روش های مبتنی بر شبکه های عصبی به طبقه بندی دسته های Emoji از Tweet بر اساس شبکه های عصبی عمیق پرداخته است.

۳. روش تحقیق

۳.۱. پیش پردازش

ابزارهای مختلف Stanford CoreNLP [1] برای تجزیه و تحلیل جامع زبانی استفاده می شود. در روش پیشنهاد، ابتدا باید متن ورودی با استفاده از این کتابخانه پیش پردازش شود. در روند اجرای روش پیشنهادی عملیاتی برای آماده سازی متن جهت اجرای روش پیشنهادی نیاز به اعمالی چون حذف ایست واژه ها، حذف حروف اضافه، حذف اعداد و نشانه گذاری داریم.



شکل ۱: روندنمای تحلیل احساسات

۳.۲. داده ها

در آنالیز احساس به غنی بودن داده های انتخاب شده برای طبقه بندی باید توجه کرد [1]. در پژوهش حاضر داده ها از سایت بانک اطلاعاتی اینترنتی فیلم ها می باشد و به طور خاص برای آنالیز احساسات انتخاب شده است.^۱ در این سامانه اطلاعات تمام فیلم ها و بررسی آن ها وجود دارد. با آنالیز احساس این بررسی ها، می توان میزان رضایت مندی کاربران از فیلم ها و محتوایشان را ارزیابی کرد.

احساسات به صورت باینری مشخص شده اند، به این معنی که بررسی هایی که فیلد rating آن ها کمتر از ۵ است، امتیاز احساسی صفر و رکورد هایی که rating آن ها بزرگتر از ۷ است، امتیاز احساسی ۱ را نشان می دهند. در این مجموعه داده هیچ فیلمی بیش از ۳۰ بررسی ندارد. ۲۵۰۰ بررسی برچسب گذاری شده در داده آموزشی مشابه هیچ یک از ۲۵۰۰ بررسی فیلم ها در داده آزمایشی نیست. علاوه بر این ۵۰۰۰۰ بررسی بدون برچسب می باشند.

۳.۳. نشانه گذاری (Tokenization): در این مرحله نظرات به کلمات مستقلی که توکن نامیده می شود، شکسته می شود [۲].

۳.۴. حذف ایست واژه ها: ایست واژه ها لغاتی هستند که علی رغم تکرار فراوان در متن، از لحاظ معنایی دارای اهمیت کمی هستند. در این فاز کلمات کم اهمیت تر و یا ایست واژه ها از متون مورد پردازش، حذف می گردند. در اغلب کاربردهای متن، حذف این کلمات نتایج پردازش را بهبود می دهد. علاوه بر این از آنجا که بیشتر کاربردهای پردازش متن با حجم عظیمی از داده ها رو به رو هستند، حذف این کلمات سبب کاهش بار محاسبات و افزایش سرعت خواهد شد.

^۱ IMDB

برای کاوش کردن مجموعه بزرگی از اسناد ضروریست که اسناد پیش پردازش شوند و اطلاعات در یک ساختار داده‌ای مناسب برای پردازش‌های بعدی ذخیره شوند. در این زمینه چندین روش وجود دارند که سعی در بهره‌گیری از ساختار نحوی و معنایی متن دارند.

در بیشتر روش‌ها، اسناد به صورت مجموعه‌ای از کلمات نمایش داده می‌شوند. بیشتر روش‌های متن‌کاوی، الگوریتم‌های کاوش را روی برچسب‌های نسبت داده شده به هر سند اعمال می‌کنند. این برچسب‌ها ممکن است کلمات کلیدی استخراج شده از سند یا فقط لیستی از کلمات در سند مورد نظر باشند. برای نشان دادن کمترین اهمیت یک کلمه در یک سند معمولاً از نمایش بردار استفاده می‌شود، برای هر کلمه یک مقدار اهمیت عددی ذخیره می‌گردد. روش‌های اصلی و مهم موجود که بر اساس این ایده هستند عبارتند از: مدل فضای بردار، مدل احتمالی و مدل منطقی. چون برخی از روش‌های متن‌کاوی که بیان می‌شوند از مدل فضای برداری استفاده می‌کنند این روش را مختصراً توضیح می‌دهیم.

۳.۵. بردار کیسه کلمات (BOW)

کیف کلمات (BOW) یک مدل در پردازش زبان‌های طبیعی است که با هدف دسته‌بندی مستندات و متون استفاده می‌شود. ایده اصلی آن، به این صورت هست که به هر کدام از کلمات یک عدد Unique نسبت می‌دهیم و Feature بدست آمده بر اساس فرکانس تکرار هر کدام از کلمات به دست خواهد آمد. به طور متعارف، مجموعه‌ای از کلمات با استفاده از فرکانس کلمه یا اهمیت به عنوان ویژگی اغلب برای طبقه‌بندی یا بازیابی سند استفاده می‌شود. یکی از مشکلات این روش، انفجار ابعاد است که با افزایش کلمات، بعد هم افزایش می‌یابد [3].

۳.۶. تکنیک TF-IDF

۳.۶.۱. جدول تکرار کلمات (TF)

روش اصلی یافتن صفحات مرتبط با یک جستجو، روش سنجش تعداد تکرار (TF) یک کلمه است. هر چه که یک کلمه در یک صفحه بیشتر تکرار شده باشد، آن صفحه ارتباط بیشتری با آن کلمه دارد بنابراین اگر کاربر کلمه‌ای را جستجو کرد، صفحاتی را نمایش خواهیم داد که آن کلمه در آن‌ها بیشتر تکرار شده باشد.

۳.۶.۲. جدول معکوس تکرار در صفحات (IDF)

به ازای هر لغت، باید فرمولی را استفاده کنیم که هر چه تکرار یک لغت در یک کتاب کمتر باشد، به آن امتیاز بیشتری بدهد (رابطه معکوس) مثلاً می‌توانیم از فرمول N/DF استفاده کنیم که N تعداد کل کتابها و DF تعداد کتابهای حاوی آن لغت است. با این فرمول هر چه یک لغت کمتر تکرار شده باشد، عدد بزرگتری تولید می‌شود.

۳.۶.۳. جدول TF-IDF

با داشتن جدول TF و IDF، می‌توانیم مرتبط بودن یک لغت با یک صفحه را با ضرب این دو در هم نمایش دهیم:

$$TF * IDF = \text{میزان ارتباط لغت با یک صفحه}$$

^۲ Term Frequency

یعنی هر چه يك لغت در يك متن بيشتر به كار رفته باشد و در ساير متن‌ها خيلي كم به كار رفته باشد، امتياز آن صفحه براي آن لغت بيشتر مي شود كه منطقي هم به نظر مي رسد.

به هر كلمه در متن يك وزن اختصاص دهيم. با اين كار، مي توانيم اهميت يك كلمه را در فرايند مهندسي ويژگي^۳ بهتر شناسايي کرده و در نهايت ويژگي‌هاي بهتري را براي تزريق به الگوريتم‌هاي بعدي مانند طبقه‌بندي يا خوشه‌بندي داشته باشيم. مقدار TF-IDF مخفف دو كلمه است: TF به معني Term Frequency يعني تعداد تکرار يك كلمه در يك متن و عبارت IDF به معني Inverse Document Frequency كه مي توان آن را به برعكس تعداد تکرار در متون ترجمه كرد. وزن دهی TF-IDF طبق فرمول (۱) محاسبه می شود.

$$f(w) = TF(w).IDF(w) = TF(w). \log \frac{N}{n(w)+1} \quad (۱)$$

۴. الگوريتم هاي دسته بندي

۴.۱. الگوريتم جنگل تصادفي^۴

الگوريتم جنگل تصادفي يكي از روش‌هاي زير مجموعه درخت هاي تصميم گيري است كه براي حل مسأله، تعداد زيادي درخت تصادفي روي زير مجموعه هايي از مجموعه داده توليد مي نمايد [4] و دقت طبقه بندي اين روش با ساخت مجموعه اي از درختان و رأي گيري بين آن ها براي به دست آوردن رده اي با بيشترين تعداد رأي، پيشرفت هاي قابل توجهي داشته است [5]. جنگل تصادفي يك روش يادگيري ماشين چند منظوره است كه قادر به انجام هر دو وظيفه رگرسيون و طبقه بندي است. يك نوع از روش يادگيري گروهی است، كه در آن گروهی از مدل‌های ضعیف تركيب می‌شوند تا يك مدل قدرتمند را شكل دهند [6]. يكي از پارامترهاي مهم و موثر در دقت دسته بندي، تعداد درختان در جنگل است. هر يك از درختان موجود در جنگل به تنهائي كارايي زيادي نداشته و در واقع قدرت تشخيص جنگل به برآيند قدرت كليي درختان وابسته است [7]. يكي از مزايای جنگل تصادفي، استفاده از داده‌هاي بزرگ با ابعاد بالا است. اين متد داراي روش‌هايي براي تعادل خطاها در مجموعه داده‌هايي است كه در آن كلاس‌ها نامتعادل هستند. جنگل تصادفي شامل نمونه‌برداري از داده ورودی با جاگزینی به نام نمونه‌برداري بوت استرپ^۵ است.

۴.۲. الگوريتم درختان اضافي^۶

الگوريتم Extra-Trees يك مجموعه اي از درخت تصميم گيري يا رگرسيون غيرمعمول را ايجاد مي كند با توجه به روش كلاسيك بالا به پايين. دو تفاوت اصلي آن، با ساير روش هاي گروهی، بر پايه دسته بندي اين است كه، گره ها را با انتخاب نقطه هاي برش به صورت تصادفي جدا مي كند و از كل نمونه يادگيري (به جای بوت استرپ) براي رشد درخت استفاده می کند [8].

^۳ Feature Engineering

^۴ Random Forest Algorithm

^۵ Bootstrap Sampling

^۶ Extremely randomized trees

۴.۳. الگوریتم ماشین بردار پشتیبان^۷

ماشین بردار پشتیبان (SVM) یک الگوریتم آموزشی است. طبقه‌بند را برای پیش‌بینی نمونه جدید از کلاس، آموزش می‌دهد. SVM عمده‌تاً مبتنی بر ایده صفحات تصمیم‌گیری است که به صراحت در مورد مرز تصمیم‌گیری صحبت می‌کند و حوزه‌ای که مرز تصمیم‌گیری بین کلاس‌ها می‌باشد را به عنوان پارامتر مورد استفاده قرار می‌دهد.

دو کاربرد کلیدی از تکنیک SVM وجود دارد که برنامه‌نویسی ریاضی و تابع کرنل هستند. SVM می‌تواند سطح جداکننده را بین داده‌های کلاس‌های متفاوت در فضای چند بعدی، بهینه کند. طبقه‌بند بردار پشتیبان^۸ (SVC) طبقه‌بندی هیبرید را جستجو می‌کند. SVC به عنوان تابع کرنل روی سطوح تصمیم‌گیری غیر خطی کاربرد دارد [9].

در واقع SVM یک Hyperplane که بین نمونه‌های مثبت و منفی مجموعه آموزش قرار می‌گیرد را مشخص می‌کند. پارامترهای b_j به گونه‌ای تنظیم می‌شوند که فاصله بین Hyperplane و نزدیک‌ترین نمونه مثبت و منفی ماکزیمم شود. که طبق فرمول (۲) محاسبه می‌گردد.

$$y = b_0 + \sum_{j=1}^n b_j t_{aj} \quad (2)$$

۴.۴. الگوریتم نایو بیز^۹

Naïve Bayes یک الگوریتم یادگیری ماشین برای حل مشکلات طبقه‌بندی است. برای ساخت مدل و پیش‌بینی سریع از قضیه احتمال استفاده می‌کند. اساس قضیه Bayes، برپایه استقلال ویژگی‌ها است. یک روش مبنا برای طبقه‌بندی متن و حل مشکل قضاوت اسناد درباره تعلق یک سند به یک دسته (مانند هرزنامه یا مشروع، ورزش و یا سیاست، و غیره) است که از روش تعیین فرکانس کلمه به عنوان ویژگی استفاده می‌کند. پیش پردازش مناسب، در این الگوریتم با روش‌های پیشرفته تر از جمله ماشین‌های بردار پشتیبان [10] در رقابت است.

۵. نتیجه و ارزیابی

در این مقاله هدف، بررسی و مقایسه اثر بخشی الگوریتم‌های داده کاوی، جهت آنالیز احساس به روش نظارت شده مبتنی بر داده‌های برچسب‌گذاری انجام شده است. در این راستا از الگوریتم‌های جنگل‌های تصادفی، درخت اضافی، بردار پشتیبان و نایو بیز استفاده شده است که با بررسی چگونگی عملکرد هر الگوریتم و پارامترهای تاثیرگذار آن مورد ارزیابی و بررسی قرار گرفت. در ابتدای این پژوهش، متون با دو رویکرد به فضای برداری انتقال داده می‌شود. که در رویکرد اول از تکنیک BOW و در رویکرد دوم از TF-IDF استفاده گردید. طبق جدول (۱)، در همه الگوریتم‌ها، با تکنیک TF-IDF نتایج بهتری حاصل گردید. در ضمن جنگل‌های تصادفی عملکرد بهتری نشان دادند. دلیل اصلی آن بکارگیری تکنیک ترکیب در این الگوریتم می‌باشد.

که با استفاده از تعدادی از الگوریتم‌های درخت تصمیم و تهیه دیتاست‌های مختلف با تکنیک جایگشتی و در نهایت رأی‌گیری اکثریت ما بین این درختان، نتیجه نهایی بدست آمد. با توجه به کاربرد و اهمیت نتایج حاصل از تحلیل احساسات در حوزه‌های

^۷ Support Vector Machine

^۸ Support Vector Classifier

^۹ Naïve Bayes

مختلف، به کارگیری مدل پیشنهادی جنگل‌های تصادفی توصیه شود. درضمن نظرات کاربران، منعکس کننده عقاید و نظرات واقعی آن‌ها در مورد محصولات و خدمات می‌باشد و از این جهت منبع ارزشمندی برای ارائه پیشنهاد به مخاطبین است.

جدول ۱: مقایسه تاثیر الگوریتم‌ها

Bow				
ناپو بیز	Linear SVC	ExtraTrees	جنگل تصادفی	
۰,۷۷۰۸۶	۰,۷۷۳۷۱	۰,۷۶۷۶۸	۰,۷۸۴۶۹	دقت (p)
۰,۷۶۱۷۱	۰,۷۷۳۶۷	۰,۷۶۲۵۴	۰,۷۸۰۸۵	صحت R
۰,۷۵۹۱۲	۰,۷۷۳۵۹	۰,۷۶۰۹۷	۰,۷۷۹۷۴	معیار F
TF-IDF				
ناپو بیز	Linear SVC	ExtraTrees	جنگل تصادفی	
۰,۷۸۹۸۴	۰,۷۹۱۵۹	۰,۷۷۸۵۰	۰,۸۰۵۴۵	دقت (p)
۰,۷۸۹۸۴	۰,۷۹۱۶۰	۰,۷۷۲۱۸	۰,۸۰۳۳۴	صحت R
۰,۷۸۹۷۹	۰,۷۹۱۵۹	۰,۷۷۰۴۲	۰,۸۰۲۶۶	معیار F

۶. منابع

- [1] C. D. Manning, J. Bauer, J. Finkel, and S. J. Bethard, "The Stanford CoreNLP Natural Language Processing Toolkit," pp. 55–60, 2014.
- [2] T. Tran, D. Nguyen, A. Nguyen, and E. Golen, "Sentiment Analysis of Emoji-based Reactions on Marijuana-Related Topical Posts on Facebook," *2018 IEEE Int. Conf. Commun.*, pp. 1–6, 2018.
- [3] K. Matsumoto, "Classification of Emoji Categories from Tweet Based on Deep Neural Networks," pp. 17–25, 2018.
- [4] E. Scornet, "A random forest guided tour," vol. 25, no. 2, pp. 197–227, 2016.
- [5] "طبقه بندی ترافیک شبکه با استفاده از الگوریتم جنگل تصادفی بهبودیافته", خویی. ا. ز. p. 15, 30 11 1395.
- [6] S. RAY, "Powerful Guide to learn Random Forest(With codes in R & Python)," *Powerful Guide to learn Random Forest(With codes in R & Python)*, p. 7, 7 SEPTEMBER 2015.
- [7] مشخصات نویسندگان مقاله بررسی تعداد درختان در الگوریتم جنگل تصادفی برای تحلیل عقاید در فروشگاههای "مسلمی. م. مجازی", p. 6, 01 01 1395.
- [8] A. Tyryshkina, N. Coraor, and A. Nekrutenko, "Predicting runtimes of bioinformatics tools based on historical data : Five years of Galaxy usage," pp. 1–10, 2019.
- [9] "Immune Support Vector Machine Approach for Credit Card(svm).pdf."
- [10] A. Varsani and K. Manoj, "Advance and Innovative Research," no. April, 201