

Adversarial Attacks on a Text Sentiment Analysis Model

Sahar Mokarrami Sefidab¹, Seyed Abolghasem Mirroshandel^{2*}, Hamidreza Ahmadifar³, Mehdi Mokarrami Sefidab⁴

1. Faculty of Engineering, University of Guilan, Rasht, Iran. sahar.mokarrami94@gmail.com
2. Assistant Professor, Faculty of Engineering, University of Guilan, Rasht, Iran. (*Corresponding Author*)
mirroshandel@guilan.ac.ir
3. Assistant Professor, Faculty of Engineering, University of Guilan, Rasht, Iran. ahmadifar@guilan.ac.ir
4. Payam-e Noor University of Guilan, Rasht, Iran. info.mokarrami@gmail.com

Abstract

Background and Purpose: Recently some researchers have shown that deep learning models, despite their high accuracy, can be vulnerable through some manipulations of their input samples. This manipulation leads to the production of new samples called Adversarial examples. These samples are very similar to the original ones, so humans cannot differentiate between these samples and the original, and cannot remove them from the dataset before predicting the model and preventing model errors. Various types of research have been done to generate malicious samples and inject them into the model, among which, the production of text samples has its own difficulties due to the discrete nature of the text. In this research, we tried to reach the highest level of vulnerability by providing a method with the least manipulation of the input data, and by testing the proposed method, we were able to bring the accuracy of CNN and LSTM models to less than 10%.

Methods: In this research, for making malicious samples, first, a word that can increase the amount of error in the classification prediction is selected from the word dictionary as a candidate word for replacement by using Taylor expansion and then considering the importance of each word in the calculated cost of the corresponding candidate word, we proposed an arrangement for substitution between words. Finally, we moved the words in the specified order until the output of the model changed.

Results: The evaluation of the presented method on two sentiment analysis models, LSTM and CNN, has shown that the proposed method has been very effective in reducing the accuracy of both models to less than 10% with a small number of replacements and this indicates the success of the proposed method compared to some other similar methods.

Conclusion: As mentioned, most of the attention of science and industry is on the production of different systems using deep learning methods, so their security of them is also important. It is important to increase the strength of the models against adversarial examples. In this research, a method with the least amount of manipulation was presented to produce textual conflict samples. It seems that in the future it will be possible to use different methods of making natural texts to produce samples that, in addition to the apparent similarity to the original sample, are also comprehensible in terms of content.

Keywords: Text Attacks, Adversarial Examples, Loss function gradient, Sentiment Analysis, Natural Language Processing.

حملات تخصصی در یک مدل تحلیل احساس متن

سحر مکرمی سفیدآب^۱، سیدابوالقاسم میرروشندل^{۲*}، حمیدرضا احمدی فر^۳، مهدی مکرمی سفیدآب^۴

۱. دانشکده فنی، دانشگاه گیلان، رشت، ایران. sahar.mokarrami94@gmail.com

۲. استادیار، دانشکده فنی، دانشگاه گیلان، رشت، ایران. (نویسنده مسئول) mirroshandel@guilan.ac.ir

۳. استادیار، دانشکده فنی، دانشگاه گیلان، رشت، ایران. ahmadifar@guilan.ac.ir

۴. دانشگاه پیام‌نور واحد رشت، رشت، ایران. info.mokarrami@gmail.com

چکیده: شبکه‌های عصبی عمیق دقت و کارایی بالایی در حل مسائل مختلف دارند، اما در برابر نمونه‌های تخصصی آسیب‌پذیر هستند. این دسته از نمونه‌های مخرب به‌منظور فریب مدل آموزش‌دیده و بررسی آسیب‌پذیری مدل‌های شبکه عصبی تولید می‌شوند. در حوزه متن روش‌های موفق زیادی برای ساخت این نمونه‌ها ارائه‌نشده‌است. در این پژوهش یک روش قوی مبتنی بر گرادیان تابع هزینه مدل برای تولید نمونه‌های تخصصی متنی ارائه شده و نشان داده‌است که می‌توان با جایگزینی تعداد کمی از کلمات موجود در نمونه‌های اصلی با کلماتی که بیشترین تأثیر منفی را روی تصمیم طبقه‌بند دارند، نمونه‌های جدیدی مشابه با نمونه‌های اولیه برای فریب طبقه‌بند تحلیل احساس در سطح کلمه تولید نمود. در نهایت با بهره‌گیری از این نمونه‌ها دقت دو مدل طبقه‌بند از پیش آموزش‌دیده، بر روی روش مورد استفاده در این پژوهش، با دست‌کاری اندک نمونه‌های ورودی، موفق به کاهش دقت طبقه‌بندی از ۸۶ درصد به کمتر از ۱۰ درصد شده‌است.

واژه‌های کلیدی: حملات متنی، نمونه‌های تخصصی، گرادیان تابع هزینه، تحلیل احساس، پردازش زبان طبیعی.

۱. مقدمه

در دنیای امروز رویکردهای یادگیری عمیق تقریباً در تمامی حوزه‌های هوش مصنوعی به خصوص تجزیه و تحلیل احساس، توانسته‌اند با ارائه مدل‌های پهنه، راه‌حل‌های سریع و دقیقی را برای حل مسائل مختلف داشته‌باشند. درحالی‌که برخی مطالعات نشان داده‌است که مدل‌های یادگیری عمیق در برابر نمونه‌های تخصصی به شدت آسیب‌پذیر هستند به طوری که با تشخیص اشتباه خسارات جبران ناپذیری را تحمیل می‌نمایند. نمونه‌های تخصصی، نوع تغییر یافته و دست‌کاری شده نمونه‌های ورودی اولیه خود هستند که باعث فریب مدل‌های شبکه عصبی و رخداد خطا در آن‌ها می‌شوند [۱]. این قبیل نمونه‌ها شباهت بسیاری زیادی با نمونه‌های اولیه خود دارند به همین دلیل انسان نمی‌تواند بین این نمونه‌ها با نمونه‌های اصلی موجود در مجموعه دادگان تفاوتی قائل شود؛ بنابراین نمی‌تواند قبل از پیش‌بینی مدل آن‌ها را از مجموعه داده‌ها حذف نماید تا مانع از خطای مدل در پیش‌بینی شود. به همین منظور باید به دنبال روشی برای مبارزه با این نمونه‌ها و مقاوم‌سازی مدل‌های شبکه عصبی در برابر آن‌ها بود.

تلاش‌های مختلف برای ساخت نمونه‌های تخصصی از یک سو و استفاده از این نمونه‌ها در مقاوم‌سازی و بررسی آسیب‌پذیری مدل‌های یادگیری عمیق از سویی دیگر منجر به گشایش در چپه جدیدی از دانش در حوزه یادگیری عمیق شده‌است. با استفاده از این نمونه‌ها می‌توان به بررسی امنیت مدل‌های شبکه عصبی پرداخته و آسیب‌پذیری مدل‌های مختلف متنی را نشان داد. از طرفی با یافتن محدودیت‌های مدل می‌توان ارزیابی و تفسیر مناسبی از آن مدل داشت. مثلاً در پژوهشی از این نمونه‌ها برای ارزیابی مدل‌های درک مطلب متنی به منظور توسعه مدل‌های دقیق‌تر [۲] و در پژوهشی دیگر برای مشاهده و سنجش میزان آسیب‌پذیری مدل‌های ترجمه ماشینی استفاده شده‌است [۳].

در زمینه تصویر پژوهش‌های بسیاری برای ساخت نمونه‌های تخصصی ارائه شده‌اند که اکثر آن‌ها موفق بوده‌اند. اما در حوزه متن این پژوهش‌ها موفقیت چندانی نداشته‌اند چرا که متن برخلاف تصویر ماهیتی گسسته دارد. هر تصویر از تعداد زیادی پیکسل تشکیل شده‌است که اعمال تغییرات کوچک روی آن‌ها و ایجاد اختلال در آن‌ها از دید انسان پنهان مانده و درعین حال می‌تواند مدل را نیز در تشخیص، همراه نماید. اما در حوزه متن جابه‌جایی و تغییر در سطح کلمه یا حرف منجر به تغییر اساسی در ساختار متن می‌شود و مفهوم متن را بهم می‌ریزد. بنابراین استفاده مستقیم از روش‌های تولید نمونه‌های تخصصی تصویری برای ایجاد اختلال و دست‌کاری داده‌های متنی کارآمد نیست [۴].

بر همین اساس روش‌های مختلفی برای ساخت نمونه‌های تخصصی و ایجاد حمله روی مدل‌های متنی ارائه شده‌اند. این روش‌ها به دو نوع کلی حملات جعبه سیاه و جعبه سفید دسته‌بندی می‌شوند. در حملات جعبه سیاه، مدل شبکه عصبی و گرادیان آن به صورت کامل در دسترس

نیست. بنابراین در این نوع از حملات از روش‌های آزمون و خطا برای تولید نمونه مخرب استفاده می‌شود به گونه‌ای که پس از تولید نمونه‌ای خروجی مدل را بررسی نموده تا زمانی که نمونه مخرب مناسب یافت شود، اما برخلاف حملات جعبه سیاه در حملات جعبه سفید مهاجم به مدل شبکه عصبی دسترسی کامل دارد بنابراین در این قبیل حملات از روش‌های مبتنی بر گرادیان به عنوان یک رویکرد حمله استفاده شده و پژوهشگران با روش‌های مختلفی تلاش می‌کنند تا داده‌ها را به گونه‌ای که گرادیان تابع هزینه افزایش یابد، دست‌کاری نمایند و در نتیجه منجر به خطای مدل در پیش‌بینی شوند [۵]. حملات تخصصی از نظر گاهی دیگر به دو صورت حملات هدفمند و حملات بی‌هدف نیز تقسیم‌بندی می‌شوند. در حملات هدفمند در جستجوی روشی هستیم تا مدل تقریباً به‌ازای بیشتر نمونه‌ها یک برچسب هدف را پیش‌بینی نماید درحالی‌که در حملات بدون هدف تنها از مدل می‌خواهیم تا در تشخیص برچسب نمونه‌ها دچار خطا شود [۶].

برقراری غلط‌های املائی روی کلمات، حذف، اضافه کردن و یا جایگزینی کلمات با کلمات مشابه دیگر، از جمله اولین کارهای انجام‌شده برای تولید نمونه‌های تخصصی در حوزه متن هستند. به عنوان نمونه، در یک حمله در سطح کاراکتر، یک متن ساده به صورت انتخاب یک کاراکتر اطراف کاراکتر مورد نظر در ساختاری مانند ساختار صفحه کلید و جایگزینی آن با کاراکتر اصلی معرفی شد. این مدل روشی برای حمله روی یک سیستم ترجمه ماشینی محسوب می‌شود. به طور مشابه در روشی دیگر، نمونه‌هایی برای فریب Google Perspective API از طریق ایجاد برخی غلط‌های املائی روی لغات ورودی یا اضافه کردن نقطه یا کاراکترهای مشابه به حروف تولید شد. در این روش‌ها اغلب اوقات نمونه‌های تولید شده از دید انسان پنهان نیست [۳،۷]. راهکار دیگری برای جایگزینی کلمات، استفاده از الگوریتم ژنتیک است که به صورت تکراری جمعیتی از راه‌حل‌های کاندید را به راه حل بهتر تکامل می‌دهد. بدین شکل که کلمه‌ای انتخاب شده و سپس با استفاده از فاصله اقلیدسی، تعدادی از نزدیک‌ترین همسایه‌های آن کلمه انتخاب می‌شود و سپس کلمه‌ای که احتمال برچسب هدف برای فریب مدل را افزایش دهد به عنوان کلمه مورد نظر برای جایگزینی انتخاب می‌شود [۸]. یک روش دیگر، وابسته به گرادیان تابع هزینه برای تولید نمونه‌های تخصصی است. در این روش ابتدا گرادیان تابع هزینه مدل به ازای تمامی کلمات موجود در متن، محاسبه می‌شود و پس از آن کلمه‌ای با بیشترین مقدار گرادیان به عنوان کلمه مناسب برای افزودن اختلال به آن و همراهی مدل انتخاب می‌شود زیرا می‌تواند خطای مدل را افزایش دهد [۹]. در روش دیگری یک استخری از کلمات مترادف، از نظر ظاهری مشابه و نیز غلط‌های املائی به‌ازای کلمات موجود، در نظر گرفته شده‌است و جایگزینی کلمات موجود در نمونه با هر کدام از کلمات مشابه استخر بر اساس اهمیت هر کلمه با محاسبه گرادیان نسبت به ورودی انجام گرفته‌است [۱۰]. در پژوهشی دیگر فضایی برای تعبیه کلمات در نظر گرفته شده و بردار اختلال کوچکی در فضای تعبیه

برداری به کلمات موجود در آن فضا افزوده شده است و سپس این بردار جدید به نزدیکترین کلمه اطراف خود در آن فضای برداری نگاشت شده است [۱۱]. در یکی از جدیدترین پژوهش‌ها برای تولید نمونه‌هایی با ساختار متنی مناسب، از یک مدل زبانی پوششی استفاده شده است. در این روش با نمونه‌گیری از توزیع احتمالاتی کلمات موجود در نمونه‌های اصلی که با کمک مدل زبانی پوششی محاسبه شده‌اند و به‌روزرسانی وزن‌های مدل زبانی مبتنی بر خروجی طبقه‌بندی، طی چند گام نمونه جدیدی مشابه با نمونه اولیه تولید شد [۴]. از سوی دیگر می‌توان برای تولید نمونه مناسب، جهت بردارهای لغاتی را که در فضای برداری قرار می‌گیرند، محدود به کلمات موجود در آن فضا نمود [۱۲].

نوع دیگری از حملات به‌عنوان حمله سرسری محسوب می‌شوند. این نوع از حملات مستقل از داده‌های ورودی هستند به‌گونه‌ای که بدون دست‌کاری و تغییر مستقیم نمونه‌های ورودی می‌توان یک اختلالی را به تمامی نمونه‌ها افزود تا به‌زای بیشتر نمونه‌ها مدل دچار اشتباه در تشخیص شود [۱۳]. برای ایجاد حملات سرسری متنی، در یکی از اولین روش‌ها، ابتدا تعدادی کلمه با طول مشخص به‌عنوان اختلال اولیه به ابتدای تمامی نمونه‌های اصلی اضافه شده، سپس برای این که اختلال خاصیت تخریب‌پذیری داشته باشد، گرادیان تابع هزینه مدل نسبت به هر کدام از کلمات اختلال به بردار کلمات افزوده می‌شود. در انتها نیز نزدیکترین بردار به بردار قبلی در فضای تعبیه برداری به‌عنوان کلمه جایگزین انتخاب می‌شود [۱۳].

در یک پژوهش، برای اینکه نمونه‌های تخصصی سراسری تولید شده ساختاری مشابه با عبارت و جمله داشته و با معنا باشند، از یک خودمرزگذار استفاده شده است. در این روش یک بردار نویز گاوسی انتخاب شده و به‌عنوان ورودی، به یک مولد داده می‌شود و مولد آن را به یک فضای پنهان منتقل می‌کند در نهایت با استفاده از رمزگشا خروجی فضای پنهان به یک عبارت با معنی تبدیل می‌شود. در هر بار تکرار این عملیات، تصویر گرادیان تابع خطای مدل حول بردار نویز، محاسبه شده و به نویز قبلی اضافه می‌شود که این کار برای به‌روزرسانی بردار نویز و افزایش خطای مدل انجام می‌گیرد [۱۴].

در روشی، برای تولید نمونه‌هایی که از نظر محتوایی مشابه با نمونه‌های اولیه هستند، ارزش هر کلمه در متن و نیز اختلاف مقادیر احتمالات برچسب نمونه اصلی با نمونه تغییر یافته، محاسبه شده و به‌صورت حریصانه، کلمات داخل متن با کلمات مترادف خود جایگزین می‌شوند [۱۵]. در روشی دیگر بهترین جایگزینی با محاسبه گرادیان تابع هزینه نسبت به آن جایگزینی و انتخاب بیشترین مقدار گرادیان، مشخص شده و در نهایت با استفاده از جستجوی پرتویی یک دسته از بهترین کلمات برای فریب مدل و جایگذاری در متن ورودی انتخاب می‌گردد [۱۶، ۱۷]. در یک روش دیگر از الگوریتم متروپلیس هستینگز برای تولید کلمات مناسب استفاده شده است. البته در برخی

حالات استفاده از این روش منجر به تولید جملات و عبارات ناقص شده است [۵، ۱۸].

در این پژوهش روشی قوی و وابسته به گرادیان تابع هزینه به‌صورت هدفمند به‌منظور ایجاد تغییراتی روی داده‌های ورودی ارائه می‌نماییم و نشان خواهیم داد که با استفاده از این روش تنها با جایگزینی تعداد بسیار کمی از کلمات موجود در نمونه‌های ورودی با کلماتی دیگر، می‌توان دقت طبقه‌بندی را به‌شدت کاهش داد.

روش ارائه شده به‌صورت جایگذاری کلمات است که در آن هزینه هر جابه‌جایی به‌همراه ارزش هر کلمه در نمونه ورودی برای ایجاد یک حمله متنی مناسب در نظر گرفته می‌شود. در واقع برای ترتیب‌دادن یک حمله قدرتمند، با استفاده از بیشترین هزینه هر جابه‌جایی بر اساس اهمیت لغات موجود در متن ورودی، توانسته‌ایم کلمات کاندید را با کلمات موجود در نمونه‌های اصلی جابه‌جا کنیم تا جایی که نمونه مورد نظر در طرف دیگر مرز تصمیم طبقه‌بندی قرار گیرد و مدل را در تشخیص دچار خطا کند. با این کار نیازی به جستجو در میان کلمات کاندید برای یافتن یک ترتیب مناسب از لغات برای جایگزینی نیست. در نتیجه این روش نسبت به برخی دیگر از روش‌ها سرعت اجرای بالایی نیز دارد.

روش پیشنهادی بدین صورت است که در ابتدا به‌زای هر کدام از کلمات موجود در نمونه ورودی، از میان مجموعه کلمات موجود در لغت نامه، کلمه‌ای که بتواند میزان خطا در پیش‌بینی طبقه‌بندی را افزایش دهد به‌عنوان کلمه کاندید برای جایگزینی انتخاب‌نموده و سپس با در نظر گرفتن اهمیت هر کلمه در مقدار هزینه محاسبه شده کلمه کاندید مربوط، یک ترتیبی برای جایگزینی میان کلمات پیشنهاد می‌دهیم. در نهایت روش ارائه شده را روی دو مدل تحلیل احساسات مختلف ارزیابی‌نموده و نشان داده‌ایم که این روش می‌تواند با دست‌کاری اندک نمونه‌های ورودی، مدل‌ها را وادار به اشتباه در تشخیص کند. آزمایش‌های مختلف روی روش ارائه شده نشان داده است که روش پیشنهادی در کاهش دقت هر دو مدل بسیار مؤثر بوده است به‌طوری که این روش توانسته دقت مدل را به کمتر از حدود ۱۰٪ کاهش دهد و این امر بیانگر موفقیت روش پیشنهادی نسبت به برخی دیگر از روش‌های مشابه است.

ساختار مقاله بدین‌صورت است که ابتدا در بخش دوم به شرح مفاهیم و تعاریف اولیه در رابطه با نمونه‌های تخصصی پرداخته و سپس روش مورد استفاده در این پژوهش را در بخش سوم بیان می‌نماییم. در نهایت در بخش چهارم به تشریح نتایج حاصل و در بخش پنجم به جمع‌بندی مطالب پژوهش خواهیم پرداخت.

۲. بیان مسئله

طبقه‌بند شبکه عصبی f می‌تواند هر ورودی $x_i \in X = \{x_1, \dots, x_n\}$ را به برچسب $y_i \in Y = \{y_1, \dots, y_k\}$ با میزان احتمال بالایی نگاشت نماید. با افزودن اختلال بسیار کوچک روی نمونه‌های ورودی و دست‌کاری داده‌های ورودی می‌توان طبقه‌بندی را مورد حمله قرار داد تا دچار خطا شده و برچسبی غیر از برچسب اصلی را پیش‌بینی نماید. با این کار

می‌توان مدل‌های شبکه عصبی را بررسی و ارزیابی نموده و آن‌ها را نسبت به حملات تخاصمی مقاوم و آگاه نمود. میزان دست‌کاری روی داده‌های ورودی باید بسیار کم باشد از طرفی این تغییر باید بتواند دقت مدل را کاهش دهد. حال هرچقدر میزان تغییرات اعمال‌شده بر روی داده‌های ورودی به‌منظور گمراهی مدل کمتر باشند و بتوانند دقت مدل را به‌طرز مطلوبی کاهش دهند، حمله بهتری خواهیم داشت [۱۶].

۱.۲. حملات هدفمند

در حملات تخاصمی بدون هدف، با اعمال تغییرات اندک روی داده‌های ورودی مدل را وادار نموده تا برچسبی غیر از برچسب اصلی را پیش‌بینی نماید؛ در صورتی که در یک حمله هدفمند با دست‌کاری داده‌های ورودی، مدل را وادار می‌نماییم تا یک برچسب مشخص و هدف، مانند y_{target} را پیش‌بینی نماید. اگر فرض کنیم نمونه تغییر یافته به‌صورت x_{adv} باشد، در این حالت در جستجوی راهکاری برای حداقل‌سازی تفاوت میان پیش‌بینی مدل روی نمونه تغییر یافته و y_{target} با استفاده از دست‌کاری داده ورودی و در نتیجه محاسبه مسئله بهینه‌سازی زیر هستیم [۱۳]:

$$\arg \min_{x_{adv}} \text{loss}(y_{target}, f(x_{adv})) \quad (1)$$

در رابطه ۱، $\text{loss}(y_{target}, f(x_{adv}))$ تابع خطای مدل را نمایش می‌دهد و این در واقع فاصله میان خروجی مدل روی نمونه تغییر یافته شده و y_{target} را محاسبه می‌کند.

۳. روش پیشنهادی

در این بخش جزئیات روش پیشنهادی خود را به‌منظور تغییر داده‌های ورودی و تولید نمونه‌های تخاصمی ارائه خواهیم کرد و پس از تولید نمونه‌های تخاصمی، با تزریق آن‌ها به مدل شبکه عصبی، آن‌را دچار خطا در تشخیص می‌کنیم. هر متن ورودی از ترکیب تعدادی کلمه تشکیل می‌شود، بنابراین برای دست‌کاری ورودی در سطح کلمه، یک روش، جایگزینی بعضی از کلمات متن با کلمات مخربی است که بتوانند منجر به فریب مدل و در نتیجه تشخیص نادرست در خروجی مدل شوند. در ادامه به بیان چگونگی یافتن کلمات و جایگزینی با کلمات موجود در نمونه می‌پردازیم.

۱.۳. یافتن کلمات کاندید

در مرحله اول می‌بایست کلماتی را به‌عنوان کاندید برای جایگزینی با هر کدام از کلمات موجود در متن اصلی بیابیم. این کلمه کاندید باید به‌گونه‌ای باشد تا بعد از جایگزینی بتواند مدل را فریب دهد. برای این کار، به‌ازای هر کلمه موجود در متن ورودی، از میان کلمات موجود در مجموعه لغت نامه، کلمه‌ای را به‌عنوان جایگزین انتخاب می‌کنیم که بیشترین مقدار هزینه را به مدل تحمیل نماید [۱۶، ۱۷].

یک روش مناسب برای یافتن کلمات جایگزینی که بیشترین تأثیر منفی روی تصمیم طبقه‌بند دارند، بهره‌گیری از بسط تیلور مرتبه اول

می‌باشد. با بهره‌گیری از بسط تیلور می‌توان سریع‌تر هزینه جابه‌جایی هر کدام از کلمات با کلمات موجود در لغت نامه را تخمین زد و به‌جای این‌که بعد از هر بار جابه‌جایی کلمات، هزینه دقیق را با درخواست‌زدن به مدل محاسبه‌نماییم، یک بار گرادینان تابع هزینه نسبت به مقادیر ورودی را طی یک فاز برگشت به عقب به‌دستی آوریم.

به عبارت دقیق‌تر در این حالت گرادینان هزینه مدل نسبت به هر یک از کلمات ورودی طی یک فاز برگشت به عقب محاسبه‌شده و بردار خروجی حاصل، در ماتریس تعبیه برداری که در آن تمامی کلمات لغت‌نامه به صورت برداری نگاشت شده‌اند، ضرب می‌شود. در نتیجه با این کار برداری خواهیم داشت که در آن هزینه جابه‌جایی هر کلمه از متن اصلی با کلمات موجود در لغت‌نامه محاسبه‌شده و سپس با انتخاب بیشترین مقدار هزینه به‌ازای هر کلمه به‌همراه اندیس متناظر با آن، برداری به‌دستی‌آید که مناسب‌ترین جایگزین را به‌ازای هر لغت برای فریب مدل در خود دارد. رابطه ۲ این امر را نشان می‌دهد.

$$\arg \max_{v_i \in V} [v_i]^T \cdot \nabla_x \text{loss} \quad (2)$$

در این رابطه V بردار تعبیه تمامی کلمات موجود در لغت‌نامه است و $\nabla_x \text{loss}$ نیز بیانگر میانگین گرادینان خطای مدل روی نمونه‌های ورودی است.

۲.۳. جایگذاری کلمات

در مسائل طبقه‌بندی متنی، هر کلمه موجود در متن، دارای ارزش و تأثیری در تصمیم‌گیری طبقه‌بند است. به‌عنوان مثال در عبارت «such a good movie» کلمه good نسبت به سایر کلمات اهمیت بیشتری دارد. در نتیجه طبقه‌بند با ارزیابی آن عبارت، کل متن را به‌صورت یک متن مثبت پیش‌بینی می‌کند. حال بعد از آن‌که به‌ازای هر یک از کلمات موجود در متن اصلی، کلمه کاندیدی برای جایگزینی یافت‌شده، در مرحله بعدی برای مشخص نمودن اولویت جابه‌جایی هر کدام از کلمات با کلمات کاندید خود، ارزش تمامی کلمات موجود در متن اصلی را به‌صورت جداگانه محاسبه‌نموده و به‌صورت یک بردار وزنی در بردار هزینه جابه‌جایی کلمه مورد نظر که در مرحله قبل محاسبه‌نمودیم، ضرب می‌کنیم. در نهایت نیز مقادیر به‌دست‌آمده را به‌ترتیب بزرگترین مقادیر مرتب می‌کنیم و با این کار به هر کلمه‌ای برای جایگزینی با کلمه متناظر خود امتیازی دهیم تا کلمات بر اساس امتیاز کسب شده خود جایگزین شوند تا زمانی که مدل دچار خطا شود.

ارزش هر کلمه را که با $S(x, w_i)$ نمایش داده می‌شود، می‌توان از طریق محاسبه اختلاف احتمال وقوع برچسب واقعی روی نمونه ورودی و احتمال وقوع برچسب واقعی در صورت حذف کلمه مذکور و قراردادن توکن $\langle \text{UNK} \rangle$ در آن جایگاه، به‌صورت رابطه زیر تعریف کرد [۱۵]:

$$S(x, w_i) = p(y_{true}|x) - p(y_{true}|\hat{x}_i) \quad (3)$$

به‌طوری‌که:

$$x = w_1, \dots, w_i, \dots, w_d$$

$$\hat{x}_i = w_1, \dots, \langle UNK \rangle, \dots, w_d$$

مدل شبکه عصبی پیچشی با دو لایه پنهان (CNN) و لایه آخر این مدل نیز یک لایه کاملاً متصل برای محاسبه کلاس نهایی استفاده شده است.

در هر دو مدل از یک لایه تعبیه برای نگاشت نمونه‌های ورودی به فضای پیوسته و بردار تعبیه کلمات word2vec با فضای تعبیه ۳۰۰ بعدی استفاده نموده ایم [۲۰].

۱.۴. ارزیابی

الگوریتم حمله پیشنهادی در این پژوهش را پس از پیاده‌سازی، با الگوریتم PWWS، به عنوان یک روش پایه مقایسه‌شده [۱۵]، رویکرد ارائه شده در PWWS یکی از روش‌هایی است که توانسته نسبت به سایر روش‌های ارائه شده در حوزه تولید نمونه‌های تخصصی، با میانگین نرخ جا به جایی و تغییرات اندک روی متن ورودی، موفق عمل نماید. بنابراین در پژوهش‌های مختلف، به عنوان یک روش پایه برای سنجش و مقایسه استفاده می‌شود. از طرفی این الگوریتم با دست‌کاری و جا به جایی نمونه‌های ورودی موفق شده به طرز مطلوبی دقت مدل را روی برخی از پیکره‌های ارائه شده برای طبقه‌بندی متون نیز کاهش دهد.

روش حمله PWWS به صورت جایگزینی کلمات به طرز حریصانه است به طوری که در آن با محاسبه اهمیت هر یک از کلمات موجود در جمله ورودی و هم‌چنین میزان تأثیر جا به جایی هر کدام از کلمات با مترادف خود آن هم از طریق محاسبه اختلاف خروجی مدل روی نمونه اصلی با خروجی مدل پس از جایگزینی کلمه مورد نظر با کلمه مترادف خود، رویکردی برای جایگزینی لغات ارائه شده است. در این روش به جای انتخاب کلمات جایگزین از فضای برداری، هر کلمه با کلمه مترادف خود از مجموعه wordnet جایگزین می‌شود و در انتها کلمه‌ای که بیشترین مقدار خطا را در پیش‌بینی طبقه‌بندی دارد، به عنوان کلمه مناسب برای جایگزینی، انتخاب می‌گردد.

برای سادگی در پیاده‌سازی این روش، موجودیت‌های نامدار را نادیده گرفته و برای مقایسه، هر دوی حملات را به صورت حمله هدفمند تعریف کرده ایم. به طوری که مدل به ازای هر نمونه‌ای برچسب مخالف با برچسب اصلی خود را پیش‌بینی نماید. نتایج حاصل از پیاده‌سازی و مقایسه دو روش در جدول‌های ۱ و ۲ قرار دارند.

جدول ۱: دقت مدل قبل و بعد از حمله

| دقت مدل با روش پیشنهادی | دقت مدل با روش PWWS | دقت مدل روی داده‌های هدف قبل از حمله | مدل | کلاس داده‌های تست |
|-------------------------|---------------------|--------------------------------------|------|-------------------|
| ۱،۶۳٪ | ۴۷،۶۶٪ | ۸۵،۲۸٪ | LSTM | منفی (۱) |
| ۰،۰۵٪ | ۴۳،۹۱٪ | ۸۸،۲۸٪ | LSTM | مثبت (۰) |
| ۱۰،۰۴٪ | ۵۱،۶۳٪ | ۷۷،۳۳٪ | CNN | منفی (۱) |
| ۷،۲٪ | ۷۷،۲۵٪ | ۸۶،۰۳٪ | CNN | مثبت (۰) |

در رابطه (۳)، x ، متن ورودی به طول d کلمه و \hat{x}_i متن ورودی است که در آن توکن $\langle UNK \rangle$ جایگزین کلمه w_i شده است. به همین صورت، ارزش همه کلمات موجود در متن، با استفاده از رابطه ۳ محاسبه می‌گردد. در نتیجه به ازای متن x یک بردار $\phi(x)$ خواهیم داشت. در ادامه، مطابق با رابطه ۴، بردار نهایی را در بردار هزینه کلماتی که در مرحله قبل محاسبه شده، ضرب نموده و خروجی نهایی را به ترتیب نزولی مرتب می‌کنیم. این ترتیب به دست آمده نشان‌دهنده ترتیب جایگذاری هر یک از کلمات کاندید با کلمه متناظر خود در متن اصلی است. در واقع با این کار کلمات را بر اساس اهمیت و بیشترین میزان خطا جایگذاری نموده ایم. در انتها نیز پس از هر جایگذاری یک کلمه با کلمه کاندید خود به ترتیب اولویت یافت شده، خروجی مدل محاسبه می‌شود و جا به جایی کلمات تا زمانی که پیش‌بینی مدل تغییر کند، ادامه می‌یابد. این رویکرد نسبت به روش‌های جستجو برای جا به جایی کلمات مناسب، سریع‌تر است. رابطه ۴ تابع محاسبه اولویت جا به جایی کلمات را نشان می‌دهد.

$$H(x, x_i^*, w_i) = \phi(S(x))_i \cdot \phi([v_i]^T \cdot \nabla_x \text{loss})_i \quad (4)$$

در این رابطه $\phi(z)_i$ بیانگر تابع سافتمکس روی هر کدام از بردار مقادیر گرادیان و ارزش کلمات است. از این تابع برای نرمال‌سازی هر کدام از مقادیر ارزش و هزینه، بین صفر و یک استفاده می‌نماییم که مطابق با رابطه (۵) محاسبه می‌شود.

$$\phi(z)_i = \frac{e^{z_i}}{\sum_{k=1}^K e^{z_k}} \quad (5)$$

۴. نتایج آزمایش‌ها

مجموعه داده گان. در این پژوهش از پیکره با نک درختی تحلیل احساس دانشگاه استنفورد استفاده موده ایم [۱۹]. این پیکره شامل مجموعه‌ای از نظرات بینندگان فیلم به زبان انگلیسی است که به صورت دو کلاسه و پنج کلاسه آن موجود است و در این آزمایش از نوع دو کلاسه این پیکره (SST-2) به صورت نظرات مثبت با برچسب ۰ و نظرات منفی با برچسب ۱ استفاده نموده ایم.

شبکه‌های قربانی. در این پژوهش از دو مدل شبکه عصبی به عنوان شبکه‌های قربانی برای اجرای الگوریتم روی آن‌ها به صورت زیر استفاده کرده ایم:

مدل حافظه طولانی کوتاه مدت با دو لایه پنهان (LSTM) که اغلب از این مدل برای مسائل طبقه‌بندی متون استفاده می‌شود. در این مدل لایه آخر یک لایه کاملاً متصل است که خروجی لایه پنهان را دریافت نموده و پیش‌بینی نهایی را محاسبه می‌نماید.

به منظور ارزیابی الگوریتم پیشنهادی و بررسی کیفیت نمونه‌های متنی تولیدشده و مقایسه دو روش حمله، در ابتدا پس از پاکسازی و پیش‌پردازش اولیه، تمامی نمونه‌های موجود در مجموعه داده‌های تست با برچسب منفی و مثبت را از یکدیگر جدا نمودیم تا بتوانیم به صورت جداگانه روی هریک الگوریتم حمله هدفمند را اجرا کنیم به گونه‌ای که مدل، برچسب هدف و مخالف با برچسب اصلی را پیش‌بینی کند (برای داده‌های منفی برچسب هدف را ۰ و برای داده‌های مثبت برچسب هدف را ۱ در نظر گرفته ایم).

در مرحله بعدی، مدل از پیش‌آموزش دیده خود را بررسی نموده و قبل از اعمال حمله، دقت مدل را به صورت جداگانه روی هریک از داده‌های با برچسب مثبت و منفی اندازه‌گیری کردیم. مطابق با نتایج جدول ۱ دقت مدل LSTM روی مجموعه داده‌های با کلاس مثبت در حدود ۸۸٪ و برای داده‌های منفی در حدود ۸۵٪ شده است. دقت مدل CNN نیز روی داده‌های با کلاس مثبت حدود ۸۶٪ و برای داده‌های منفی حدود ۷۷٪ گزارش شده است. حال در انتها پس از دست‌کاری و تغییر نمونه‌های موجود در مجموعه داده‌ها و اعمال الگوریتم حمله روی هر کدام از مدل‌ها می‌بایست با محاسبه دقت مدل روی داده‌های تغییر یافته میزان گمراهی مدل‌هایمان را مشخص نماییم. بعد از اجرای الگوریتم حمله، هر چقدر دقت مدل در تشخیص درست کمتر باشد، حمله قوی‌تری خواهیم داشت. با مقایسه دو روش، دیده می‌شود که روش پیشنهادی در کاهش دقت بسیار خوب عمل کرده است و توانسته دقت مدل LSTM را روی هر دو نوع داده مثبت و منفی، به زیر ۲٪ برساند و برای مدل CNN نیز دقت به زیر ۱۰٪ رسیده است. با توجه به نتایج آمده در جدول ۱، دقت مدل LSTM روی داده‌های منفی با روش پیشنهادی به ۱٫۶۳٪ و روی داده‌های مثبت ۰٫۰۵٪ رسیده است در حالی که دقت این مدل با روش PWWS روی دو داده مثبت و منفی به حدود ۴۰٪ رسیده است همین‌طور روی مدل CNN نیز دقت روی داده‌های منفی پس از حمله به حدود ۱۰٪ و برای داده‌های مثبت به حدود ۷٪ رسیده است که در مقایسه با الگوریتم PWWS به شدت افت دقت داشته‌ایم که این امر، نشان‌دهنده مؤثر بودن روش حمله پیشنهادی است.

برای این که حمله تخصصی مناسب باشد، باید تغییرات اعمال شده روی نمونه‌های ورودی بسیار اندک باشند تا از دید انسان این میزان تغییرات قابل تشخیص نباشد. روش پیشنهادی در این پژوهش توانسته با جابه‌جایی و تغییر تعداد کمی از کلمات نیز در گمراهی هریک از مدل‌ها و کاهش دقت، بسیار موفق عمل نماید.

جدول ۲ میانگین نرخ جایگزینی کلمات پس از اجرای حمله هدفمند روی هر کدام از داده‌های با برچسب مثبت و منفی را نشان می‌دهد. میزان نرخ جابه‌جایی کلمات برابر با تعداد جایگزینی‌های کلمات در یک نمونه بر کل تعداد کلمات موجود در نمونه مورد نظر است. هر چقدر میانگین نرخ جابه‌جایی کلمات با اعمال الگوریتم حمله پایین‌تر باشد، بهتر است چرا که به نوعی نامحسوس بودن این تغییرات از دید

انسان را نشان می‌دهد. در روش پیشنهادی میانگین جابه‌جایی‌ها برای فریب مدل LSTM روی داده‌های منفی و مثبت به ترتیب ۱۳٫۹۸٪ و ۱۰٫۸۹٪ شده است، در حالی که در روش PWWS این میزان حدود ۱۹٪ و ۱۷٪ گزارش شده است و در مدل CNN نیز میانگین نرخ جابه‌جایی روی داده‌های منفی و مثبت به ترتیب ۹٫۴۵٪ و ۸٫۱۵٪ شده است، در حالی که در روش PWWS این میزان حدود ۲۵٪ و ۱۷٪ گزارش شده است. روش پیشنهادی در این مورد نیز روی هر دو نوع مدل شبکه عصبی نسبت به روش PWWS بهتر عمل کرده است.

جدول ۲: میانگین نرخ جابه‌جایی کلمات پس از حمله

| میانگین نرخ جابه‌جایی کلمات در روش PWWS | میانگین نرخ | میانگین نرخ جابه‌جایی کلمات در روش PWWS | مدل | کلاس |
|---|-------------|---|----------|-----------------|
| ۱۳٫۹۸٪ | ۱۹٫۲۲٪ | LSTM | منفی (۱) | مجموعه داده تست |
| ۱۰٫۸۹٪ | ۱۷٫۸۴٪ | LSTM | مثبت (۰) | |
| ۹٫۴۵٪ | ۱۷٫۸۳٪ | CNN | منفی (۱) | |
| ۸٫۱۵٪ | ۲۵٫۴۸٪ | CNN | مثبت (۰) | |

برای مقایسه نمونه‌های تولیدی و تغییر یافته با نمونه‌های اصلی و تصویرسازی بهتر، در جدول ۳، تعدادی از نمونه‌های تولیدشده با استفاده از روش پیشنهادی و روش PWWS را نشان داده‌ایم و خروجی مورد نظر را روی مدل LSTM اندازه‌گیری نموده‌ایم. به ازای هر کدام از حملات هدفمند و دست‌کاری نمونه‌های موجود در داده‌های منفی و مثبت، یک نمونه اصلی را همراه نمونه تغییر یافته‌اش قرار داده‌ایم. در هر یک از مثال‌ها پیش‌بینی مدل LSTM، قبل و بعد از اجرای الگوریتم حمله و نیز میزان درجه اطمینان برای هر پیش‌بینی نیز اندازه‌گیری شده است. کلمات تغییر یافته در جدول با رنگ قرمز مشخص شده‌اند.

از بررسی و مقایسه موارد موجود در جدول می‌توان مشاهده کرد که میزان تغییرات صورت گرفته روی نمونه‌های اصلی با روش پیشنهادی نسبت به روش PWWS بسیار کمتر است. مثلاً در روش PWWS نمونه اول جدول که توسط مدل به صورت منفی پیش‌بینی شده، با چهار تغییر از کلاس منفی به کلاس مثبت انتقال یافته و با درجه اطمینان ۶۰٫۶۸٪ مدل را گمراه کرده است. در حالی که در روش پیشنهادی، همان نمونه تنها با یک تغییر توانسته با درجه اطمینان نسبتاً بالایی (۹۲٫۵۵٪) مرز تصمیم مدل LSTM را جابه‌جا کرده و منجر به تشخیص خطا شود. برای یک نمونه مثبت نیز می‌توان مشاهده کرد که در روش PWWS، مدل با درجه اطمینان کمتر (۵۳٫۲٪) و تغییرات بیشتری (سه تغییر) دچار خطا شده است. در حالی که روش پیشنهادی تنها با یک تغییر و درجه اطمینان بیشتر از ۹۵٪ موفق به گمراهی مدل شده است.

جدول ۳: مثالی از نمونه‌های تولیدشده

| نمونه تغییر یافته | نمونه اصلی | پیش بینی بعد از حمله | پیش بینی قبل از حمله | روش حمله |
|---|--|---------------------------------------|---------------------------------------|-----------------|
| zaidan's script has just decent plot to string the stunts together and not rather enough word to keep the faces straight. | zaidan's script has barely enough plot to string the stunts together and not quite enough characterization to keep the faces straight. | مثبت (۰) درجه اطمینان ٪۶۰٫۸۶ | منفی (۱) درجه اطمینان ٪۹۸٫۳۱ | روش PWWS |
| zaidan's script has barely enough plot to string the stunts together and infectiously quite enough characterization to keep the faces straight. | zaidan's script has barely enough plot to string the stunts together and not quite enough characterization to keep the faces straight. | مثبت (۰) درجه اطمینان ٪۹۲٫۵۵ | منفی (۱) درجه اطمینان ٪۹۸٫۳۱ | روش پیشنهادی |
| Seldom has a movie therefore nearly matched the tone of a man and his work. | Seldom has a movie so closely matched the spirit of a man and his work. | منفی (۱) درجه اطمینان ٪۵۳٫۲۰ | مثبت (۰) درجه اطمینان ٪۹۶٫۲۰ | روش PWWS |
| ill-constructed has a movie so closely matched the spirit of a man and his work. | Seldom has a movie so closely matched the spirit of a man and his work. | منفی (۱) درجه اطمینان ٪۹۵٫۴۰ | مثبت (۰) درجه اطمینان ٪۹۶٫۲۰ | روش پیشنهادی |

استفاده از گرادیان تابع هزینه، از میان لغات موجود در لغت نامه انتخاب- نموده و با در نظر گرفتن بیشترین اهمیت هر لغت همراه میزان هزینه محاسبه شده برای جایگذاری آن‌ها اولیتی برای جایگزینی کلمات ارائه- شد.

نتایج آزمایش‌ها نشان داده‌اند که با اجرای الگوریتم پیشنهادی روی نمونه‌های ورودی و دست‌کاری نمونه‌ها و تزریق نمونه‌های تغییر یافته به مدل، هریک از مدل‌ها دچار افت دقت شدیدی شدند طوری که این میزان برای مدل LSTM به کمتر از ۲٪ و برای یک مدل CNN به کمتر از ۱۰٪ رسیده‌است. از سوی دیگر، با محاسبه میانگین نرخ جابجایی کلمات نشان داده شد که با استفاده از روش پیشنهادی می‌توان تنها با جایگزینی تعداد کمی از کلمات موجود در متن ورودی با کلماتی دیگر، مدل را به اشتباه‌انداخته و وادار به تشخیص نادرست نمود.

در آینده تلاش خواهیم کرد تمرکز بیشتری روی محتوای متون تولیدی داشته‌باشیم و کلمات موجود در متن را با کلماتی متناسب و مرتبط با متن ورودی جایگزین نماییم طوری که علاوه بر پایین نگه‌داشتن دقت مدل و اعمال تغییرات اندک، متون تولید شده از لحاظ مفهومی نیز قابل تشخیص توسط انسان را نباشند و با ارائه روش‌های مختلف دفاعی بتوانیم مدل‌های شبکه عصبی را نسبت به این قبیل حملات آگاه و مقاوم‌سازیم. از سویی دیگر، نمونه‌های تخصصی تولید شده علاوه بر بحث مقاوم‌سازی مزایای جانبی زیادی نیز دارند یکی از این موارد تقویت دادگان می‌باشد، بنابراین به نظر می‌رسد با بهبود روش‌های تولید نمونه‌های تخصصی در آینده بتوانیم از این قبیل نمونه‌ها در جهت تقویت دادگان ناکافی نیز استفاده نماییم.

مراجع

- [1] C.Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhun, I. Goodfellow and R. Fergus, "Intriguing properties of neural networks", 2nd International Conference on Learning Representations, ICLR 2014, Banff, Canada, 2014.
- [2] R. Jia., P. Liang, Adversarial examples for evaluating reading comprehension systems. In EMNLP, 2017
- [3] Y. Belinkov and Y. Bisk, "Synthetic and natural noise both break neural machine translation", In Proceedings of ICLR, 2018.
- [4] I. Fursov, A. Zaytsev, P. Burnyshev, E. Dmitrieva, N. Klyuchnikov, A. Kravchenko, E. aArtemova and E. Burnaev, "A differentiable language model adversarial attack on text classifiers", arXiv:2107.11275v1 [cs.CL], 23 Jul 2021.
- [5] Z. Kong, J. Xue, Y. Wang, L. Huang, Z. Niu and E. Li, "A survey on adversarial attack in the age of artificial intelligence", Wireless Communications and Mobile Computing, Volume 2021, Article ID 4907754, 22 pages, 2021.
- [6] J. Xu and Q. Du, "TextTricker: Loss-based and gradient-based adversarial attacks on text classification models", Engineering Applications of Artificial Intelligence, Volume 92, Elsevier, 0952-1976, 2020.

با توجه به نتایج آزمایش‌ها می‌توان مشاهده کرد که روش پیشنهادی در تمام موارد بررسی شده از جمله کاهش میزان دقت مدل، ایجاد تغییرات بسیار کم و تغییر مرز تصمیم طبقه‌بند با درجه اطمینان بالا، توانسته بسیار موفق عمل نماید. به طوری که مقایسه‌های مختلف با روش PWWS بیانگر مؤثر بودن این روش در فریب مدل‌ها و یافتن نقاط ضعف و آسیب‌پذیری‌های شبکه عصبی نسبت به این قبیل از نمونه‌ها است. بنابراین می‌توان از این رویکرد برای تولید نمونه‌های متنی، مشابه با نمونه‌های اولیه و با کمترین میزان تغییر، برای بررسی آسیب‌پذیری‌های مدل و در انتها مقاوم سازی مدل‌های شبکه عصبی نسبت به این قبیل از نمونه‌ها استفاده کرد.

۵. نتیجه‌گیری و کارهای آتی

به دلیل گسسته بودن ماهیت متن، ساخت نمونه‌های تخصصی متنی امری چالش‌برانگیز است. تغییرات اعمال شده روی نمونه‌های متنی نیز باید اندک باشند تا از دید انسان پنهان مانده و به سادگی قابل شناسایی و تشخیص نباشند. در این پژوهش، روشی بهینه برای تولید نمونه‌های تخصصی متنی به منظور حمله روی مدل‌های یادگیری عمیق ارائه شد. که در آن برای یافتن کلمات جایگزین، به‌ازای هر کدام از کلمات موجود در متن اصلی، کلمه‌ای که بیشترین هزینه را در هر جایگزینی دارد با

- [20] T. Mikolov, E. Grave, P. Bojanowski, C. Puhersch and A. Joulin, "Advances in pre-training distributed word representations", In LREC, 2018.
- [7] H. Hosseini, S. Kannan, B. Zhang and R. Poovendran, "Deceiving google's perspective api built for detecting toxic comments," arXiv preprint arXiv:1702.08138, 2017.
- [8] M. Alzantot, Y. Sharma, A. Elgohary, B. Ho, M. Srivastava and K. Chang, "Generating natural language adversarial examples, in Proceedings of Conference on Empirical Methods in Natural Language Processing (EMNLP), 2018.
- [9] B. Liang, H. Li, M. Su, P. Bian, X. Li and W. ChangShi, "Deep text classification can be fooled", arXiv preprint arXiv:1704.08006, 2017.
- [10] S. Samanta and S. Mehta, "Towards crafting text adversarial samples", arXiv preprint arXiv:2003.10388, July 2017.
- [11] N. Papernot, P. McDaniel, A. Swami and R. Harang, "Crafting adversarial input sequences for recurrent neural networks", In 2016 IEEE Military Communications Conference, MILCOM 2016, Baltimore, MD, USA, p.p. 49–54, November 1-3, 2016.
- [12] M. Sato, J. Suzuki, H. Shindo and Y. Matsumoto, "Interpretable adversarial perturbation in input embedding space for text", In Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, (IJCAI 2018), Stockholm, Sweden, p.p. 4323–4330, July 13-19, 2018.
- [13] M. Behjati, S. M. Moosavi-Dezfooli, M. SoleymaniBaghshah and P. Frossard, "Universal adversarial attacks on text classifiers", In ICASSP, 2019.
- [14] L. Song, X. Yu, H. Peng and K. Narasimhan, "Universal adversarial attacks with natural triggers for text classification", arXiv:2005.00174v2 [cs.CL], 7 Apr 2021.
- [15] S. Ren, Y. Deng, H. He and W. Che, "Generating natural language adversarial examples through probability weighted word saliency", In Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, Florence, Italy, p.p. 1085–1097, 2019.
- [16] J. Ebrahimi, A. Rao, D. Lowd and D. Dou, "Hotflip: White-box adversarial examples for text classification", In Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics, Volume 2: Short Papers, p.p. 31–36, 2018.
- [17] E. Wallace, S. Feng, N. Kandpal, M. Gardner and S. Singh, "Universal adversarial triggers for attacking and analyzing nlp", arXiv preprint arXiv:1908.07125, 2019.
- [18] H. Zhang, H. Zhou, N. Miao and L. Li, "Generating fluent adversarial examples for natural languages", in Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, Florence, Italy, 2019.
- [19] R. Socher, A. Perelygin, J. Wu, J. Chuang, C. D. Manning, A. Ng and C. Potts, "Recursive deep models for semantic compositionality over a sentiment treebank", in Proceedings of the conference on empirical methods in natural language processing (EMNLP), p.p. 1631–1642, 2013.