

# شناسایی نظرات سمی با استفاده از یادگیری انتقالی\*

محمد لشکری

## ۱ چکیده

با گسترش شبکه‌های اجتماعی در میام مردم، یکی از دغدغه‌های اصلی امروز در دنیا تبدیل این فضا به مکانی آرام و ایمن برای تبادل اطلاعات و به اشتراک گذاری نظرات است. یکی از عوامل برهم زدن نظم و آرامش فضای مجازی، منتشر کردن پیام‌های توهین‌آمیز، نژادپرستانه، تهدید و... است. ما این گونه پیام‌ها را سمی می‌نامیم. نیاز امروز ما شناسایی این نظرات سمی و جلوگیری از انتشار آن‌هاست؛ از آنجا که مدل‌های هوش مصنوعی کلاسیک دقت خوبی در این مسئله ندارند، به کارگیری مدل‌های یادگیری ماشین پردازش زبان‌های طبیعی برای حل این مسئله لازم است. در این مقاله، مدل‌هایی مبتنی بر BERT را با استفاده از یادگیری انتقالی، برای یافتن نظرات سمی/غیر سمی روی دادگان مسئله تنظیم<sup>۱</sup> کردیم. برای تنظیم، هیچ یک از لایه‌های BERT ثابت فرض نشد و همه لایه‌ها روی دادگان مسئله تنظیم شدند. مدل‌هایی که ساخته و بررسی شدند از نظر مفهومی به دو دسته کلی تقسیم می‌شوند: ۱. یادگیری توالی<sup>۲</sup> ۲. استخراج ویژگی<sup>۳</sup>. یادگیری توالی مبتنی بر شبکه عصبی تکرارشونده LSTM و استخراج ویژگی به کمک شبکه عصبی پیچشی<sup>۴</sup> یک بُعدی انجام شد. ما دو مدل مبتنی بر یادگیری توالی پیشنهاد می‌دهیم که علاوه بر LSTM از یک لایه CNN نیز استفاده می‌کند تا هم یادگیری توالی و هم استخراج ویژگی را داشته باشیم. دو مدل دیگر پیشنهادی ما فقط مبتنی بر مفهوم استخراج ویژگی ساخته شده‌اند و تنها از CNN استفاده می‌کنند. بهترین مدل به دست آمده BERT+2CNN-1D نام دارد که دو لایه CNN یک بُعدی به BERT اضافه می‌کند. این مدل، مصرف حافظه معقول و روی دادگان تست دقت بالایی دارد.

## ۲ مقدمه

مردم روزانه از شبکه‌های اجتماعی برای به اشتراک گذاری نظرات، تبادل اطلاعات، یافتن دوست و... استفاده می‌کنند. اما به اشتراک گذاری برخی نظرات که آن‌ها را **نظرات سمی** می‌نامیم، شبکه‌های اجتماعی را به بستری برای فعالیت‌های تنفرآمیز تبدیل می‌کند. نظرات سمی شامل نظرات توهین‌آمیز، غیر محترمانه و یا غیر منطقی می‌شوند که کوچک‌ترین پیامد ناشی از به اشتراک گذاری آن‌ها در شبکه‌های اجتماعی، غم‌انگیز کردن کاربران است. از این رو در

\*Toxic comment detection using transfer learning

<sup>1</sup>Fine-tune

<sup>2</sup>Sequential learning

<sup>3</sup>Feature extraction

<sup>4</sup>Convolutional neural network (CNN)

طی چند سال گذشته، محققین زیادی روی موضوع تشخیص این نوع نظرات و جلوگیری از انتشار آن‌ها کار کرده‌اند. کوتاهی جملات، مهم‌بودن ترتیب کلمات، شناسایی ارتباطات بین کلمات در فضای پنهان از جمله چالش‌هایی هستند که باعث می‌شوند روش‌های مبتنی بر فراوانی کلمات (مانند TF-IDF) عملکرد خوبی برای حل این مسئله نداشته باشند [۲].

سند  $D$  ورودی مسئله و  $C = \{0, 1\}$  مجموعه برچسب‌هاست که ۰ نشان‌دهنده سمی نبودن و ۱ نشان‌دهنده سمی بودن یک نظر است. مدل،  $c \in C$  را به سند  $D$  تخصیص می‌دهد. مدل زبانی BERT<sup>۵</sup> می‌تواند انتخاب خوبی برای حل این مسئله باشد. زیرا دادگانی که BERT روی آن از پیش آموزش دیده است شامل برخی متن‌ها و عبارات سمی بوده است [۵]. برای یافتن نظرات سمی، روش‌های متفاوتی با استفاده از پردازش زبان‌های طبیعی، یادگیری ژرف، گراف‌های دانش<sup>۶</sup> ارائه شده است. پرکاربردترین شبکه‌های عصبی برای شناسایی این گونه نظرات CNN و LSTM هستند که به ترتیب به استخراج ویژگی (در اینجا شناسایی عبارات توهین‌آمیز) و یافتن ارتباطات کلمات می‌پردازند. اگرچه هنوز در کشف نظرات سمی عملکرد عمومی‌سازی خوبی ندارند [۶]. در مقاله [۶] استراتژی‌های مختلف تنظیم BERT<sup>۷</sup> روی دادگان مسئله مطرح شده است. یکی از مدل‌های پیشنهادی آنان اضافه کردن یک لایه LSTM به مدل BERT بوده است که طبق نتایج گزارش شده در مقاله، عملکرد خوبی داشته است (نام آن را BERT+LSTM گذاشتند). مدل پیشنهادی دیگری که در [۶] مطرح شده است، خروجی تمام لایه‌های BERT را دریافت می‌کند و به یک لایه CNN دو بُعدی می‌دهد که بهترین مدل معرفی شده در این مقاله است (نام آن را BERT+CNN گذاشتند). ما در این مقاله چهار استراتژی برای تنظیم BERT پیشنهاد می‌دهیم. اولین مدل، خروجی BERT را به یک لایه LSTM-CNN می‌دهد و مدل دوم آن را به یک لایه CNN-LSTM می‌دهد؛ در هر دو مدل از CNN یک بُعدی استفاده شده است؛ این دو مدل جایگزین‌های BERT+LSTM هستند. دو مدل دیگر پیشنهادی ما، به ترتیب یک و دو لایه CNN یک بُعدی به BERT اضافه می‌کنند و ورودی آن‌ها خروجی لایه آخر BERT است؛ در مدل دوم، برای لایه ترکیب<sup>۸</sup> یک هسته<sup>۹</sup> در نظر گرفته می‌شود تا بعد از یک لایه شبکه عصبی پیچشی یک بُعدی، همچنان یک ماتریس داشته باشیم؛ این دو مدل جایگزین‌های BERT+CNN هستند. دلیل پیشنهاد دو مدل اول این است که برخی پیام‌های سمی نیاز به کشف ارتباطات کلمات (پیام‌های نژاد پرستانه) و برخی دیگر فقط نیاز به استخراج ویژگی (پیام‌های سمی معمول) دارند. بنابراین ما با پیشنهاد این دو مدل سعی داریم به مدلی برسیم که هر دو نوع را شناسایی کند. دو مدل بعدی به دلیل پیچیدگی محاسبه کمتر آن‌ها (از نظر حافظه) نسبت به مدل BERT+CNN پیشنهاد شده است.

## ۳ کارهای گذشتگان

در [۱] از شبکه‌های عصبی پیچشی برای تعیین نظرات سمی استفاده شد که تعبیه کلمات<sup>۱۰</sup> مبتنی بر word2vec بوده است. دقت کار آن‌ها روی دادگان تست حدود ۰/۹۱ گزارش شده است. مدل دیگری مبتنی بر LSTM با تعبیه کلمات SpaCy ارائه شده است که دقت آن حدود ۰/۹۵ است. اما دادگان استفاده شده، منشتر نشده است و چون از

<sup>۵</sup>Bidirectional Encoder Representations from Transformers

<sup>۶</sup>Knowledge graphs

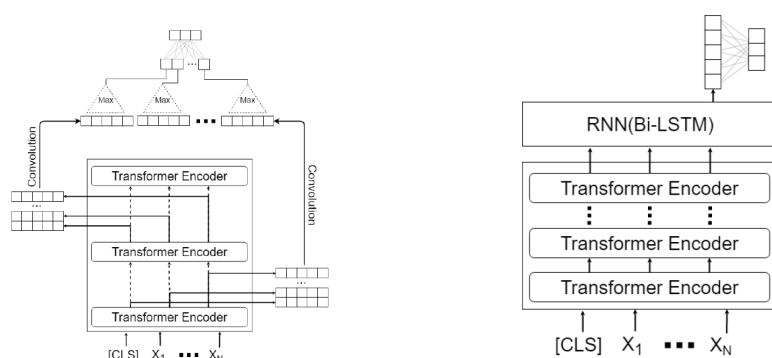
<sup>۷</sup>fine-tune

<sup>۸</sup>Pooling layer

<sup>۹</sup>Kernel

<sup>۱۰</sup>Word Embedding

یادگیری انتقالی و مدل‌های از پیش آموزش دیده در آن استفاده نشده است، نمی‌توان به عملکرد خوب عمومی‌سازی مدل مطمئن بود. در [۵] از شبکه‌های CNN و LSTM استفاده شده که تعبیه کلمات آن مبتنی بر FastText و BERT بوده است. دقت هر دو مدل در حدود ۰/۹۱ بوده و دادگان آن‌ها از پیکره توییت‌ر استخراج شده است. مدل دیگری با اضافه کردن یک لایه LSTM به انتهای BERT پیشنهاد شده است (شکل ۱ قسمت (آ)) [۶]. دقت این مدل ۰/۹۲ گزارش و پایگاه داده‌ی مورد استفاده آن‌ها از توییت‌ر جمع‌آوری شده است. مدل دیگری مبتنی بر CNN (شکل ۱ قسمت (ب)) برای تنظیم BERT ارائه شده است [۶] که دقت آن حدود ۰/۹۱ گزارش شده است. اما این مدل، پیچیدگی محاسباتی بالایی دارد و در مرحله پیش‌بینی، حافظه زیادی مصرف می‌کند.



(ب) استراتژی تنظیم BERT با اضافه کردن یک لایه CNN (BERT+CNN)

(آ) استراتژی تنظیم BERT با اضافه کردن یک لایه LSTM (BERT+LSTM)

شکل ۱: معماری‌های پیشنهادی [۶] برای تنظیم BERT

## ۴ روش پیشنهادی

مدل زبانی BERT در سال ۲۰۱۸ توسط گوگل معرفی شد که یک مدل از پیش آموزش‌دیده روی BookCorpus و Wikipedia Corpora است. از BERT می‌توان در طبقه‌بندی متن، سیستم‌های پرسش و پاسخ و کشف روابط معنایی کلمات استفاده کرد. همانطور که در بخش مقدمه گفته شد، مدل‌هایی مبتنی بر BERT با اضافه کردن یک لایه LSTM و CNN<sup>۱۱</sup> به آن برای حل این مسئله پیشنهاد شده است که نام آن‌ها به ترتیب BERT + LSTM و BERT + CNN است [۶]. ما چهار استراتژی برای تنظیم BERT روی دادگان مسئله پیشنهاد می‌دهیم:

۱. اضافه کردن یک لایه LSTM-CNN به BERT که لایه پیچشی<sup>۱۲</sup> آن یک بُعدی است و جایگزین مدل BERT + LSTM است. نام آن را BERT + LSTM-CNN می‌گذاریم.

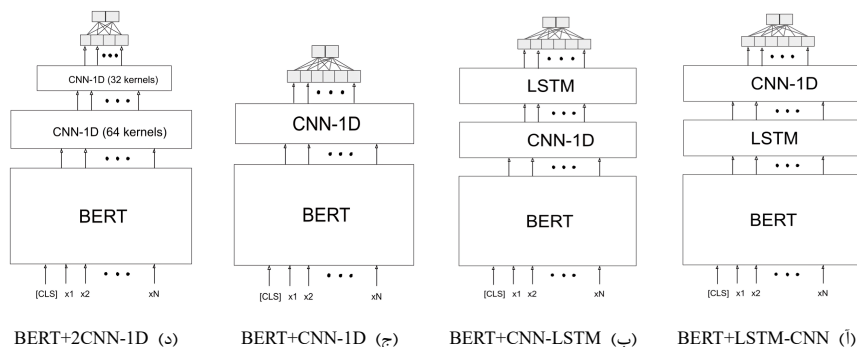
<sup>۱۱</sup> ما در این مقاله CNN را شبکه عصبی پیچشی ترجمه می‌کنیم. این شبکه از دو لایه تشکیل شده است: ۱. لایه پیچشی ۲. لایه ترکیب

<sup>۱۲</sup> Convolutional layer

۲. اضافه کردن یک لایه CNN-LSTM به انتهای BERT که لایه پیچشی آن یک بُعدی است. در این مدل ابتدا استخراج ویژگی و سپس یادگیری توالی انجام می‌شود. این مدل جایگزین BERT+LSTM است و آن را BERT+CNN-LSTM می‌نامیم.

۳. اضافه کردن یک لایه CNN به BERT که ورودی آن تنها خروجی لایه آخر BERT است و از یک لایه پیچشی و یک لایه ترکیب یک بُعدی تشکیل شده است. این مدل جایگزین مدل BERT+CNN است و آن را BERT+CNN-1D می‌نامیم.

۴. مدل پیشنهادی بعدی یک لایه CNN یک بُعدی دیگر به مدل قبل اضافه می‌کند. اگر اعمال کردن لایه ترکیب به جای کل توالی خروجی لایه پیچشی، روی هسته‌هایی از آن اعمال شود، بعد از یک لایه CNN یک بُعدی، همچنان یک ماتریس خواهیم داشت و می‌توانیم روی آن یک لایه CNN یک بُعدی دیگر اعمال کنیم. این مدل را BERT+2CNN-1D می‌نامیم.



شکل ۲: مدل‌های پیشنهادی

در شکل ۴ در مدل‌های (آ) و (ب) و (ج) تعداد هسته‌های لایه پیچشی ۶۴ است. در مدل‌های ۱ و ۲، بُعد ورودی و خروجی LSTM برابر است.

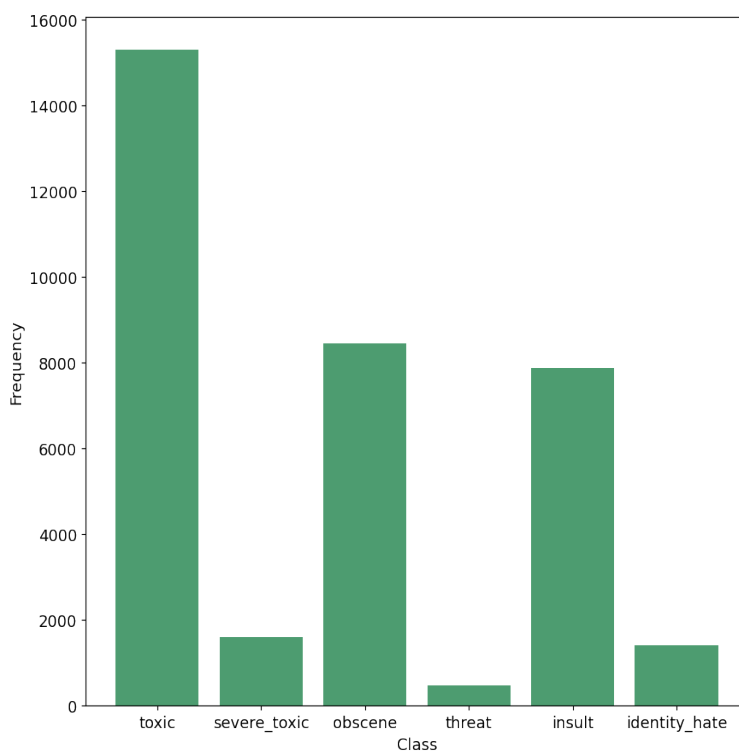
## ۵ بخش آزمایشی

### ۱.۵ دادگان

دادگان مسابقه شناسایی نظرات سمی، برگزار شده توسط گوگل در Kaggle را در نظر گرفته‌ایم که شامل ۱۵۹,۵۷۱ متن نظرات جمع‌آوری شده از توئیتر با میانگین ۶۷ کلمه برای جملات است [۴]. برچسب‌های نمونه‌ها به ترتیب toxic, severe\_toxic, obscene, threat, insult و identity\_hate هستند. اسنادی که یکی از ویژگی‌های قبل برای آنها ۱ بوده است را با برچسب ۱ (سمی) و اسنادی که هیچ یک از ویژگی‌های ذکر شده برای آنها ۱ نبوده با برچسب ۰ (غیر سمی) علامت‌گذاری می‌کنیم و مسئله را به یک مسئله طبقه‌بندی دودویی<sup>۱۳</sup> تبدیل می‌کنیم. در

<sup>13</sup>Binary class classification

اشکال ۳ و ۴ فراوانی کلاس‌ها قبل و بعد از تبدیل مسئله به طبقه‌بندی دودویی قابل مشاهده است. چون بعد از تبدیل مسئله به حالت دودویی، نامتوازن بودن شدیدی در دادگان مشاهده می‌شود، با استفاده از روش **under-sampling** دادگان را متوازن می‌کنیم. جزئیات این قسمت در بخش ۳.۵ تشریح می‌گردد.



شکل ۳: فراوانی کلاس‌ها قبل از تبدیل به طبقه‌بندی دودویی

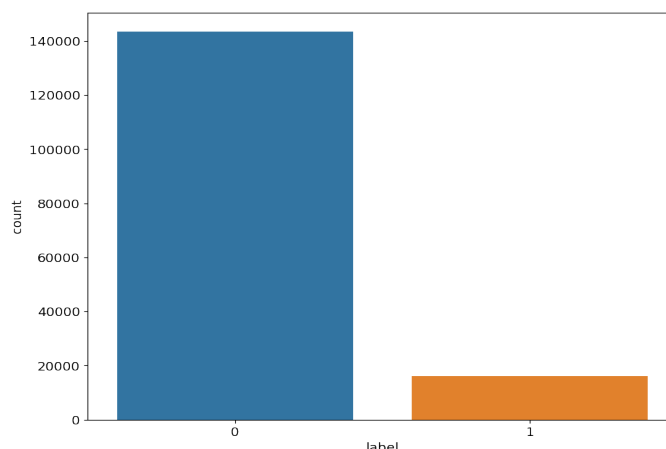
## ۲.۵ پیش‌پردازش دادگان

کاراکترهای تمام پیام‌ها به حروف کوچک تبدیل شدند. تمام اسامی کاربران، اعداد، هشتک‌ها، رایانشانی‌ها<sup>۱۴</sup> و علامت‌های احساسات به ترتیب با `<user>` و `<number>` و `<hashtag>` و `<url>` و `<emoticon>` جایگزین شدند. تمام کاراکترهای اضافی و `uni-code` ها از متون حذف شدند. چون مدل BERT توالی کلمات را یاد می‌گیرد، کلمات ایست<sup>۱۵</sup> حذف نشدند و از ساده‌سازی<sup>۱۶</sup> کلمات نیز صرف نظر شد.

<sup>۱۴</sup>URLs

<sup>۱۵</sup>Stop words

<sup>۱۶</sup>Stemming



شکل ۴: فراوانی کلاس‌ها بعد از تبدیل به طبقه‌بندی دودویی

## ۳.۵ پیاده‌سازی و نتایج

در قسمت پیاده‌سازی، برای استفاده از مدل BERT، از کتابخانه transformers و برای شبکه‌های LSTM و CNN از کلاس‌های آماده pyTorch استفاده کردیم. از محیط Google colab به عنوان بستر پیاده‌سازی استفاده شد. توکن‌سازی دادگان به وسیله کلاس BertTokenizer از کتابخانه transformers با حداکثر طول جمله ۳۶ انجام شد. دادگان با نسبت ۸۰ و ۱۰ و ۱۰ درصد به آموزش، اعتبارسنجی و تست تقسیم شدند. اندازه هسته<sup>۱۷</sup> برای لایه‌های پیچشی و ترکیب ۳ و اندازه گام<sup>۱۸</sup> ۱ انتخاب شد. وزن‌های هر لایه با نسبت ۰/۱ حذف شدند. مدل در ۳ مرحله<sup>۱۹</sup> با اندازه دسته ۳۲ به وسیله الگوریتم AdamW<sup>۲۰</sup> با نرخ یادگیری  $10^{-5} \times 2$  آموزش داده شد. وزن تمام لایه‌های BERT در تمام لایه‌ها روی دادگان مسئله تنظیم شد. جهت متوازن‌سازی دادگان، به تعداد ۲ برابر تعداد نمونه‌های کلاس یک از کلاس صفر نمونه گرفته شد و سایر پیام‌ها نادیده گرفته شد. با انجام این کار، فراوانی کلاس یک ۲۲۵، ۱۶ و فراوانی کلاس صفر ۴۵۰، ۳۲ به دست آمد. جزئیات پیاده‌سازی مدل‌ها به شرح زیر است:

۱. **مدل BERT+LSTM-CNN:** از آنجا که حداکثر طول جمله را ۳۶ قرار دادیم؛ خروجی BERT به صورت (۳۶, ۷۶۸) است. این ماتریس به یک لایه LSTM با بُعد خروجی ۳۶ داده می‌شود و بعد از ترانزاده کردن، به یک شبکه پیچشی با تعداد ۶۴ هسته داده می‌شود. خروجی نهایی یک ماتریس با بُعد (۳۴, ۶۴) است که بعد از مسطح‌سازی به یک لایه چگال داده می‌شود.

۲. **مدل BERT+CNN-LSTM:** پیاده‌سازی این مدل مشابه مورد اول است با این تفاوت که ابتدا یک شبکه پیچشی و بعد LSTM اعمال می‌شود. خروجی نهایی، یک ماتریس با بُعد (۳۴, ۶۴) است که بعد از

<sup>17</sup>Kernel size

<sup>18</sup>Stride

<sup>19</sup>Epochs

<sup>20</sup>Adam with decoupled weight decay

مسطح‌سازی به یک لایه چگال داده می‌شود.

۳. **مدل BERT+CNN-1D**: خروجی BERT به یک شبکه پیچشی با ۶۴ هسته داده می‌شود؛ بعد خروجی نهایی (۶۴, ۳۴) است. تفاوت این مدل با دو مدل قبل در این است که یادگیری توالی در آن انجام نمی‌شود و فقط یک لایه CNN جهت استخراج ویژگی به کار گرفته می‌شود.

۴. **مدل BERT+2CNN-1D**: از آنجا که لایه ترکیب را به جای اعمال روی کل توالی ورودی، روی هسته‌هایی با اندازه ۳ از آن با اندازه گام ۱ اعمال می‌کنیم، بعد از یک شبکه پیچشی یک بُعدی با ۶۴ هسته، همچنان یک ماتریس خواهیم داشت که بُعد آن (۶۴, ۳۴) است. این ماتریس به یک شبکه پیچشی دیگر با ۳۲ هسته پیچشی و لایه ترکیب با اندازه ۳ و طول گام ۱ داده می‌شود. خروجی یک ماتریس (۳۲, ۳۲) است که بعد از مسطح‌سازی به یک لایه چگال داده می‌شود.

از آنجا که ارزش هر دو کلاس برابر است، از دو معیار Accuracy و F1 استفاده می‌کنیم که F1 میانگین ماکزیمی امتیاز F1 برای هر دو کلاس است. در کنار این دو معیار از دو معیار Precision و Recall به عنوان معیارهای ثانوی استفاده می‌کنیم. فرض کنیم  $c \in \{0, 1\}$  و  $T_c$  و  $F_c$  به ترتیب، نشان‌دهنده تعداد اعضای از مجموعه آزمایشی هستند که به درستی و به اشتباه در کلاس  $c$  پیش‌بینی شده‌اند و  $F'_c$  تعداد اعضای باشد که برچسب مشاهده شده آن‌ها  $c$  بوده اما در کلاس دیگری پیش‌بینی شده‌اند.  $P_c$  و  $R_c$  و  $F1_c$  که به ترتیب صحت<sup>۲۱</sup>، دقت<sup>۲۲</sup> و امتیاز F1 کلاس  $c$  هستند که به صورت زیر محاسبه می‌شوند:

$$P_c = \frac{T_c}{T_c + F_c}, R_c = \frac{T_c}{T_c + F'_c}, F1_c = 2 \cdot \frac{P_c R_c}{P_c + R_c}$$

با توجه به فرمول‌های فوق، معیارهای ارزیابی به صورت زیر به دست می‌آیند:

$$\begin{aligned} \text{Accuracy} &= \frac{T_0 + T_1}{T_0 + T_1 + F_0 + F_1}, \\ F1 &= \frac{F1_0 + F1_1}{2}, \\ \text{Precision} &= \frac{P_0 + P_1}{2}, \\ \text{Recall} &= \frac{R_0 + R_1}{2} \end{aligned}$$

نتایج حاصل شده، در جدول زیر قابل مشاهده است. عملکرد BERT+LSTM-CNN و BERT+CNN-LSTM تقریباً به اندازه BERT+LSTM بوده است. یکی از دلایل این اتفاق می‌تواند این باشد که برای تشخیص سمی بودن یا نبودن جملات موجود در این دادگان نیاز به یادگیری توالی<sup>۲۳</sup> نداریم و با استخراج ویژگی این وظیفه به خوبی قابل انجام است. زیرا فراوانی پیام‌های نژادپرستانه و تهدید که عمدتاً نیاز به یادگیری توالی به کمک شبکه‌های عصبی تکرار شونده دارند، در این پایگاه داده کم است. مدل‌های بعدی که برای رفع مشکل پیچیدگی حافظه پیشنهاد شده اند تا اندازه خوبی این مشکل را رفع کرده‌اند، در حالی که کماکان دقت آن‌ها به خوبی مدل پیشنهادی [۶] است.

<sup>21</sup>Precision

<sup>22</sup>Recall

<sup>23</sup>Sequential learning

مدل BERT+CNN برای طبقه‌بندی حدود ۴,۶۰۰ جمله مجموعه تست، نزدیک به ۶ گیگابایت حافظه نیاز دارد. اما مدل BERT+2CNN-1D با مصرف کمتر از ۲ گیگابایت حافظه، این وظیفه را انجام می‌دهد. به دلیل مصرف حافظه معقول این مدل و دقت بالای آن، تصمیم گرفتیم یک بار دیگر آن را با حداکثر طول جمله ۶۴ آموزش دهیم؛ نتایج ارزیابی این مدل در سطر آخر جدول ۱ آمده است. در این حالت، مصرف حافظه این مدل حدود ۳ گیگابایت بوده که همچنان از نظر پیچیدگی مصرف حافظه، از مدل پیشنهادی [۶] بهینه‌تر است.

Methods	Accuracy	F1	Precision	Recall
BERT + LSTM [6]	0.91	0.90	0.90	0.91
BERT + CNN-LSTM	0.92	0.91	0.91	0.91
BERT + LSTM-CNN	<b>0.92</b>	<b>0.91</b>	<b>0.92</b>	0.90
BERT + CNN [6]	0.92	0.91	0.90	0.91
BERT + CNN-1D	0.91	0.90	0.91	0.90
BERT + 2CNN-1D	0.92	0.91	0.91	0.91
BERT + 2CNN-1D (max len.=64)	<b>0.94</b>	<b>0.93</b>	<b>0.92</b>	<b>0.93</b>

Table 1: The results

## مراجع

- [1] Georgakopoulos, Spiros V., et al. "Convolutional neural networks for toxic comment classification." Proceedings of the 10th hellenic conference on artificial intelligence. 2018.
- [2] Quan, Xiaojun, et al. "Short text similarity based on probabilistic topics." Knowledge and information systems 25.3 (2010): 473-491.
- [3] Dubey, Krishna, et al. "Toxic Comment Detection using LSTM." 2020 Third International Conference on Advances in Electronics, Computers and Communications (ICAEECC). IEEE, 2020.
- [4] <https://www.kaggle.com/datasets/julian3833/jigsaw-toxic-comment-classification-challenge?select=train.csv> [Accessed: 4-23-2022]
- [5] d'Sa, Ashwin Geet, Irina Illina, and Dominique Fohr. "Bert and fasttext embeddings for automatic detection of toxic speech." 2020 International Multi-Conference on: "Organization of Knowledge and Advanced Technologies"(OCTA). IEEE, 2020.



- [6] Mozafari, Marzieh, Reza Farahbakhsh, and Noel Crespi. "A BERT-based transfer learning approach for hate speech detection in online social media." International Conference on Complex Networks and Their Applications. Springer, Cham, 2019.