



گروه مهندسی کامپیوتر
دانشگاه صنعتی همدان

دسته بندی اسناد فارسی با استفاده از یادگیری عمیق

نام دانشجو: رامین عقیقی

نام استاد راهنما: دکتر حسن بشیری

گروه مهندسی کامپیوتر، دانشگاه صنعتی همدان
ایمیل دانشجو: ramin_aghighi@yahoo.com

چکیده

دسته بندی داده های متنی در کاربردهای مختلف بسیاری مورد توجه قرار گرفته و پژوهشگران زیادی بر روی آن کار کردند. در ابتدا از روش های سنتی مانند نایو بیز، ماشین بردار پشتیبان و نزدیک ترین همسایه استفاده می شد که بخاطر وابستگی به استخراج ویژگی ها و مشکل زمان بر بودن و پرهزینه بودن استخراج ویژگی، از سال ۲۰۱۰ اکثر پژوهشگران و مقالات به سمت روش های یادگیری عمیق رفتند که نیاز به مهندسی ویژگی ندارد. روش Text CNN تعداد مرجع زیادی را در مدل یادگیری عمیق دارد که در آن یک شبکه ی عصبی کانولوشن (CNN) برای حل مشکل طبقه بندی متن برای اولین بار معرفی شده است. حالا در این پایان نامه، قصد پیاده سازی الگوریتم های به روز یادگیری عمیق مانند CNN، RNN، LSTM و GRU را بر روی داده های همشهری داریم که در زبان فارسی از داده های معتبر به حساب می آید.

واژه های کلیدی: دسته بندی اسناد، شبکه های عصبی عمیق، شبکه های عصبی کانولوشن، همشهری

راهبردهای پیشنهادی

در پژوهش حاضر، توجه به روش های مدرن از قبیل BERT و GPT، همچنین روش های مبتنی بر گراف مورد بررسی قرار نگرفته است. با این وجود، این روش ها از اهمیت بسزایی برخوردارند و می توانند در بهبود دقت دسته بندی اسناد به زبان فارسی تأثیرگذار باشند. به همین دلیل، در ادامه تحقیقات و به منظور ارتقاء کارایی مدل ها، پیشنهاد می شود که به روش های زیر توجه شود:

- 1) بررسی مدل های پیش آموزش دیده: ارزیابی و بررسی مدل های پیش آموزش دیده نظیر BERT و GPT در زمینه دسته بندی اسناد به زبان فارسی. این مدل ها به عنوان یکی از جدیدترین رویکردها در حوزه یادگیری عمیق و پردازش زبان طبیعی شناخته شده اند.
- 2) استفاده از روش های مبتنی بر گراف: بررسی و اعمال روش های مبتنی بر گراف برای دسته بندی اسناد. این رویکرد می تواند ویژگی ها و ارتباطات پیچیده تری را در متون فارسی مدل کند.
- 3) تحقیق در زمینه انتقال یادگیری: بررسی امکان انتقال یادگیری از مدل های آموزش دیده بر روی داده های انگلیسی به مدل های دسته بندی اسناد فارسی.
- 4) افزایش حجم داده های آموزشی: تلاش برای جمع آوری و افزایش حجم داده های آموزشی به منظور افزایش تعمیم پذیری مدل ها.
- 5) تجزیه و تحلیل خطاها: تحلیل دقیق خطاها و نقاط ضعف مدل ها به منظور بهینه سازی و اصلاح راهبردها و پارامترها.

با اعمال این راهبردها، می توان به بهبود کارایی و دقت مدل ها در دسته بندی اسناد به زبان فارسی دست یافت و در جهت افزایش تنوع و کارایی مدل ها قدم برداشت.

منابع

- [1] Q. Li *et al.*, "A survey on text classification: From traditional to deep learning," *ACM Transactions on Intelligent Systems and Technology (TIST)*, vol. 13, no. 2, pp. 1-41, 2022.
- [2] K. Kowsari, K. Jafari Meimandi, M. Heidarysafa, S. Mendu, L. Barnes, and D. Brown, "Text classification algorithms: A survey," *Information*, vol. 10, no. 4, p. 150, 2019.

تایید استاد راهنما

نام و امضا استاد راهنما:

تایید تحصیلات تکمیلی:

تایید امور پژوهشی:

مقدمه

در دهه های گذشته، مدل های سنتی نظیر NB، KNN، SVM و Random Forest (RF) به عنوان دسته بندی متن مورد استفاده قرار گرفتند. از سوی دیگر، مدل های مبتنی بر یادگیری عمیق مانند Text CNN و LSTM نیز به دلیل قابلیت بهتر در عملکرد، جایگاه خود را برتر یافته اند.

در زمینه دسته بندی اسناد به زبان فارسی، با وجود پیشرفت های گسترده در حوزه زبان شناسی و پردازش زبان فارسی، تعداد محدودی از تحقیقات در این زمینه انجام شده است. این تحقیقات نشان می دهند که چالش های خاصی وجود دارد که از جمله آن ها تفاوت های فرهنگی، اصطلاحات خاص به زبان فارسی و کمبود داده های آموزشی متنوع را می توان نام برد. به منظور حل این چالش ها، تصمیم گرفته ایم از مدل هایی که در زبان انگلیسی پاسخ های قابل قبولی ارائه داده اند، بهره مند شویم و دقت این مدل ها را بر روی داده های فارسی اندازه گیری نماییم. این رویکرد نوآورانه امکان استفاده از تجربیات موفق در زبان انگلیسی را برای دسته بندی دقیق تر اسناد به زبان فارسی ایجاد می کند.

در این پژوهش، هدف اصلی ما بررسی و پیاده سازی الگوریتم های یادگیری عمیق از جمله CNN، LSTM، RNN و GRU بر روی داده های همشهری فارسی است. با اجرای این الگوریتم ها، تلاش می شود تا بهبودهای قابل توجهی در دقت و کارایی دسته بندی اسناد به زبان فارسی حاصل گردد.

اهداف و روش پژوهش

طبقه بندی متن در تجزیه و تحلیل احساسات (SA)، برچسب گذاری موضوع (TL)، طبقه بندی اخبار (NC)، پاسخگویی به سؤال (QA) و غیره مورد توجه بسیاری قرار گرفته است. همینطور مدل های یادگیری عمیق نیز در دهه اخیر برای طبقه بندی متن نسبت به سایر مدل های یادگیری ماشین مانند SVM و جنگل های تصادفی و KNN، به شدت مورد استقبال قرار گرفته است. البته با وجود کارهای بسیار در زبان های دیگر، متأسفانه در حوزه زبان فارسی کارهای کمی صورت گرفته است که همین نشان از اهمیت پیاده سازی متدهای روز شبکه های عصبی مانند CNN، LSTM و GRU در زبان فارسی می باشد که بر روی دیتاست همشهری با ده هزار سند مورد بررسی قرار خواهد گرفت و در ادامه دقت آن با سایر روش ها مورد قیاس قرار خواهد گرفت.

یافته های پژوهش

در پیاده سازی الگوریتم های یادگیری عمیق بر روی داده های همشهری فارسی، نتایج حاصل از تحقیق نشان داد که از میان الگوریتم های مورد بررسی، روش های LSTM و CNN توانستند دقت مطلوبی را در دسته بندی اسناد به زبان فارسی به دست آورند. با اعمال این دو روش، دقت حاصل از مدل بر روی داده های آزمایشی به اندازه قابل توجهی ارتقاء یافت. این نتایج نشان دهنده قابلیت بالای روش های LSTM و CNN در حل مسائل دسته بندی اسناد در زبان فارسی می باشد. این پیشرفت ها در دستیابی به دقت مطلوب، اساساً بهبود عملکرد مدل در تفکیک و دسته بندی متون فارسی را نشان می دهند و این امر می تواند اساسی در توسعه راهکارهای یادگیری عمیق برای زبان فارسی باشد.