



دانشگاه صنعتی خواجه نصیرالدین طوسی
دانشکده مهندسی برق - گروه مهندسی کنترل

به نام خدا

دانشگاه تهران - دانشکده صنعتی خواجه نصیرالدین طوسی تهران

دانشکده مهندسی برق و کامپیوتر



تشخیص اندیشه

نام و نام خانوادگی	محمد جواد احمدی
شماره دانشجویی	۴۰۱۰۰۰۸۶

فهرست مطالب

۳	۱ پاسخ پرسش دوم	
۳	۱.۱ پاسخ قسمت ۱	
۷	۲.۱ پاسخ قسمت ۲	
۱۳	۳.۱ پاسخ قسمت ۳	
۳۲	۴.۱ پاسخ قسمت ۴	

فهرست تصاویر

۴ مقایسه ساختاری RNN و LSTM	۱
۵ فرآیند تبدیل کلمه به بردار	۲
۶ لایه Embedding	۳
۶ Word2Vec	۴
۷ Glove	۵
۱۳ نمایش از مدل اول	۶
۱۸ نمودار دقت و اتلاف - مدل اول - $h=25$	۷
۱۹ نمودار دقت و اتلاف - مدل اول - $h=50$	۸
۱۹ نمودار دقت و اتلاف - مدل اول - $h=75$	۹
۱۹ نمودار دقت و اتلاف - مدل اول - $h=100$	۱۰
۲۰ ماتریس درهم‌ریختگی - مدل اول - $h=25$	۱۱
۲۰ ماتریس درهم‌ریختگی - مدل اول - $h=50$	۱۲
۲۱ ماتریس درهم‌ریختگی - مدل اول - $h=75$	۱۳
۲۱ ماتریس درهم‌ریختگی - مدل اول - $h=100$	۱۴
۲۲ نمایش از مدل دوم	۱۵
۲۹ نتیجه تابع دقت و اتلاف در حالت اعتبارسنجی با داده‌های تست (منطبق بر مقاله)	۱۶
۳۰ ماتریس درهم‌ریختگی در حالت اعتبارسنجی با داده‌های تست (منطبق بر مقاله)	۱۷
۳۱ نتیجه تابع دقت و اتلاف در حالت اعتبارسنجی با داده‌های تست (با ترم تنظیم)	۱۸
۳۲ ماتریس درهم‌ریختگی در حالت اعتبارسنجی با داده‌های تست (با ترم تنظیم)	۱۹
۳۳ نتیجه تابع دقت و اتلاف در حالت داده‌های اعتبارسنجی جداگانه	۲۰
۳۴ ماتریس درهم‌ریختگی در حالت داده‌های اعتبارسنجی جداگانه	۲۱
۳۵ نتیجه تابع دقت و اتلاف در حالت داده‌های اعتبارسنجی جداگانه	۲۲
۳۶ ماتریس درهم‌ریختگی در حالت داده‌های اعتبارسنجی جداگانه	۲۳
۳۶ نمایش از مدل پاسخگو	۲۴
۴۰ نمایش از مدل Prototype Responder	۲۵

پرسش ۲. تشخیص اندیشه

۱ پاسخ پرسش دوم

توضیح پوشه کدهای تشخیص اندیشه

کدهای مربوط به این قسمت، علاوه بر پوشه محلی کدها در این لینک گوگل کولب آورده شده است.

۱.۱ پاسخ قسمت ۱

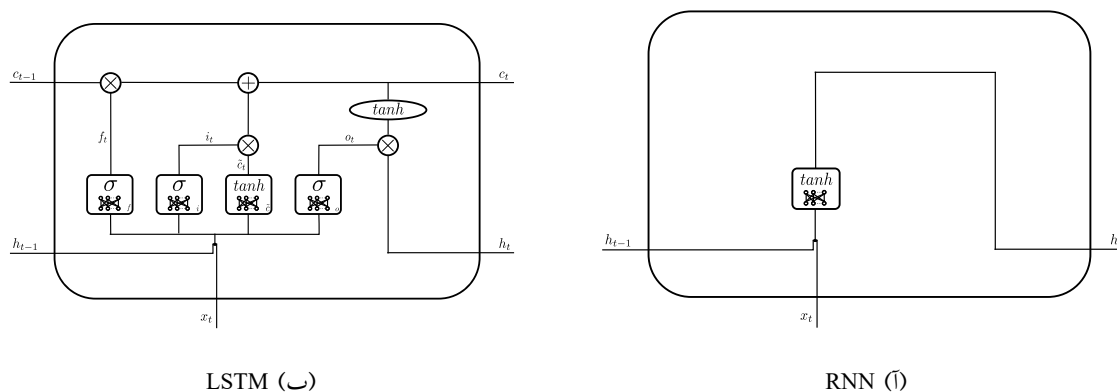
ابتدا توضیحاتی از مقاله ذکر می‌کنیم. مقاله این‌طور عنوان می‌کند که درک زبان طبیعی مستقیماً به اندیشه انسان مرتبط است و به همین دلیل ماندگار است. در حین خواندن، درک یک متن، به‌سادگی و از طریق درک یک کلمه تنها رخ نمی‌دهد، بلکه اغلب از طریق نحوه چینش کلمات اتفاق می‌افتد. به عبارت دیگر، به مدل‌سازی پویایی نیاز است که با آن کلمات فردی به وجود می‌آیند. متأسفانه، شبکه‌های عصبی پیش‌خور سنتی به‌طور مستقیم قابلیت استخراج اطلاعات از ترتیب زمانی که ورودی‌ها رخ می‌دهند را ندارند، زیرا تنها به در نظر گرفتن ورودی فعلی برای پردازش خود محدود هستند. ایده در نظر گرفتن بلوک‌های کلمات ورودی به صورت مستقل از یکدیگر در زمینه پردازش زبان طبیعی خیلی محدود است. در واقع، همان‌طور که فرآیند خواندن برای انسان منجر به حفظ تمام کلمات متن نمی‌شود، بلکه به استخراج مفاهیم اساسی بیان شده می‌انجامد، ما به یک «حافظه» نیاز داریم که فقط به در نظر گرفتن صریح ورودی‌های قبلی محدود نباشد، بلکه تمام اطلاعات مرتبط کسب‌شده در هر مرحله را به متغیرهای وضعیت تزریق می‌کند.

مجموعه این دلایل باعث به وجود آمدن شبکه‌های عصبی بازگشتی (RNN) شده است که با معرفی مفهوم «وضعیت داخلی» (بر اساس ورودی‌های قبلی)، به ویژه در وظایف پردازش زبان طبیعی، قابلیت بسیاری را نشان داده است. شبکه‌های عصبی بازگشتی حلقه‌های داخلی دارند که اجازه می‌دهند اطلاعات قبلی از طریق مراحل مختلف تحلیلی عبور کنند. می‌توان با بازکردن فعالیت‌های شبکه و در نظر گرفتن زیرشبکه‌های فردی به عنوان نسخه‌هایی (در هر مرحله به‌روز شده) از همان شبکه، که به هم متصل شده‌اند و یک زنجیره را تشکیل می‌دهند، یک دید صریح از سیستم بازگشتی به دست آورد. در این زنجیره هر یک اطلاعاتی را به جانشین خود ارسال می‌کنند.

در بحث‌های نظری، شبکه‌های عصبی بازگشتی باید قادر باشند اطلاعاتی را که الگوریتم یادگیری در یک دنباله آموزشی پیدا می‌کند، در وضعیت‌های خود نگه‌دارند. متأسفانه، در انتشار گرادیان به عقب در طول زمان، تأثیر مقادیر خروجی مطلوب که به تابع هزینه کمک می‌کنند، به حدی کوچک می‌شود که پس از تنها چند مرحله، چالش و مسأله محو گرادیان‌ها پدیدار می‌شود. این مشکل، در برخی پژوهش‌ها به صورت مشخص مورد بررسی قرار گرفته است، به این معنی است که شبکه‌های عصبی بازگشتی تنها می‌توانند یک حافظه کوتاه‌مدت داشته باشند، زیرا بخش‌های دنباله که در زمان دورتر قرار دارند، به تدریج کم‌اهمیت‌تر می‌شوند. این باعث می‌شود که شبکه‌های عصبی بازگشتی تنها برای دنباله‌های بسیار کوتاه مفید باشند که این مطلوب نیست.

برای غلبه بر مشکل حافظه کوتاه‌مدت (حداقل به طور جزئی)، ساختارهای با حافظه بلندمدت (LSTM) معرفی شدند. برخلاف یک شبکه RNN کلاسیک (تصویر ۲۳ه)، LSTM یک ساختار سلولی بسیار پیچیده‌تر دارد (تصویر ۲۳و) که در

آن یک شبکه عصبی تکی با چهار شبکه جایگزین شده که با یکدیگر تعامل دارند. با این حال، عنصر ممتاز LSTM سلول حالت (Cell State) c است که به اطلاعات اجازه می‌دهد از طریق عملیات خطی ساده در طول زنجیره جریان یابند. اضافه یا حذف کردن اطلاعات از حالت توسط سه ساختار، به نام «دروازه‌ها»، که هر کدام هدف‌های خاصی دارند، تنظیم می‌شود. اولین



شکل ۱: مقایسه ساختاری RNN و LSTM.

دروازه که «دروازه فراموشی» نام دارد، برای تصمیم‌گیری درباره اطلاعاتی که باید از حالت حذف شود، استفاده می‌شود. برای دستیابی به این هدف، حالت به طور نقطه‌ای با مقدار زیر ضرب می‌شود:

$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f), \quad (1)$$

این مقدار از یک بلوک خطی (دروازه affine به دلیل وجود انحرافات و بایاس‌ها) به دست می‌آید که بردار مشترک ورودی را با خروجی قبلی $[h_{t-1}, x_t]$ ترکیب می‌کند، و سپس تابع فعال‌سازی سیگموئیدی معمولی $(\sigma(x) = \frac{1}{1+e^{-x}})$ را را وارد داستان می‌کند. دروازه فراموشی امکان حذف (ارزش‌های نزدیک به صفر) یا حفظ (ارزش‌های نزدیک به یک) مؤلفه‌های هر بردار حالت را فراهم می‌کند. دروازه دوم که دروازه ورودی نام دارد، هدفی دارد که شرایط اضافه کردن اطلاعات جدید به حالت را فراهم کند. این عملیات از طریق ضرب نقطه‌ای بین دو بردار به دست می‌آید:

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i), \quad (2)$$

$$\tilde{c}_t = \tanh(W_c \cdot [h_{t-1}, x_t] + b_c)$$

بردار اول (همواره از طریق یک فعال‌سازی سیگموئیدی به دست می‌آید) که برای به‌روزرسانی مقادیر تصمیم می‌گیرد، و بردار دوم (از طریق یک لایه با فعال‌سازی تانژانت هایپربولیک به دست می‌آید) که هدف آن ایجاد کاندیدای‌های جدید است. توجه کنید که تابع تانژانت هایپربولیک نیز هدفی دارد که آن، تنظیم جریان اطلاعات و اجبار فعال‌سازی‌ها برای باقی ماندن در بازه منفی یک و یک است. باید توجه داشت که به‌روزرسانی وضعیت سلول تنها به دو دروازه‌ای که تازه تعریف شده‌اند وابسته است و در واقع توسط رابطه ۳ نمایش داده می‌شود:

$$c_t = f_t * c_{t-1} + i_t * \tilde{c}_t \quad (3)$$

در نهایت، «دروازه خروجی» کنترل تولید خروجی جدید h_t برای این سلول را به‌صورت وابسته به ورودی کنونی و حالت قابل مشاهده قبلی کنترل می‌کند:

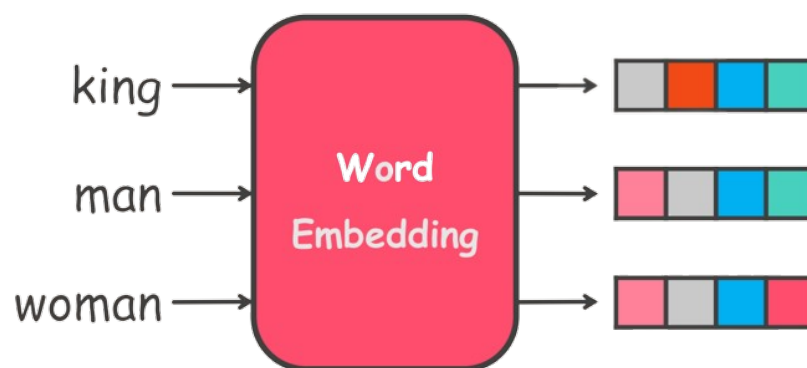
$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o), \quad (4)$$

$$h_t = o_t * \tanh(c_t).$$

هم چنین، در این حالت نیز یک لایه فعال سازی سیگموئیدی وجود دارد که قسمتی از حالت را مشخص می کند و آن ها را در مقادیر بین صفر و یک ضرب می کند. توجه کنید که، به منظور محدود کردن مقادیر خروجی، تابع تانژانت هایپربولیک به هر عنصر از بردار حالت اعمال می شود (این یک تابع ساده بدون هیچ لایه عصبی دیگری است)؛ پیش از ضرب در بردار مشخص شده توسط دروازه.

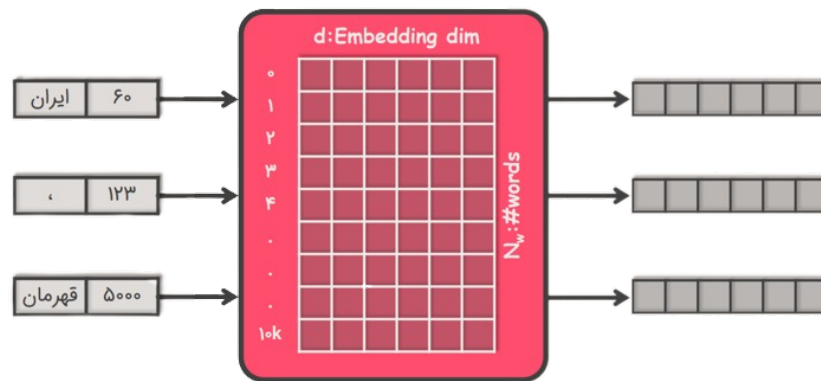
بنابراین با توجه به موارد ذکر شده، اولین مشکل RNN مسأله محو گردایی است که یکی از متهم های اصلی آن همان تانژانت هایپربولیک به کار رفته در قسمت فعال سازی است. این مسأله منجر به ضعف در حفظ وابستگی های بلند مدت می گردد. مقاله در مورد روش های کدگذاری این گونه بیان می کند که بدون شک و بدون وابستگی به نوع شبکه عصبی استفاده شده، لازم است تعیین کنیم که چه نوع تعبیری از جملات استفاده می شود. در واقع، برخلاف زبان های رسمی که به طور کامل مشخص هستند، زبان طبیعی از نیازهای ساده ارتباطی بشر نشأت می گیرد و به همین دلیل دارای تعداد زیادی ابهام است. برای حداقل سعی در فهمیدن آن، لازم است نوعی «نزدیکی معنایی» بین اصطلاحات مختلف مشخص شود، مثل تبدیل کردن هر کلمه به برداری با مقادیر واقعی و به صورت مناسب. بنابراین، تعبیر به دست آمده می تواند تک تک کلمات را به یک نمایش عددی تبدیل کند که معنای آن ها را حفظ می کند و در بعضی معانی، آن را قابل فهم برای رایانه نیز می کند. امروزه روش های مختلفی برای به دست آوردن این فضای معنایی وجود دارد که به طور کلی به عنوان تکنیک های تعبیر کلمه شناخته می شوند. در این مقاله، تصمیم گرفته شده که از یک تعبیر پیش آموزش دیده به نام GloVe استفاده شود، که بر اساس مجموعه واژگان ۴۰۰۰۰۰ کلمه ای در یک فضای ۳۰۰ بعدی نگاشت شده است.

تبدیل کلمات به بردار از اهمیت بالایی برخوردار است. ما می خواهیم کلمات را همانند آن چه که در ؟؟ نشان داده شده به یک لایه Word Embedding بدهیم و از آن طرف بردار تحویل بگیریم. این بردارهای خروجی نباید One-Hot باشند؛ چرا که، مثلاً در حالت عادی ارتباطاتی بین کلمات وجود دارد. مثلاً بین کلمات مرد و زن و یا خواهر و برادر ارتباطاتی است. دو بردار One-Hot همواره ضرب نقطه ای و تشابه صفر دارند و به درد ما نمی خورند. بنابراین خروجی بردار باید چیزی مثل همان خروجی برداری در شکل ۲ باشد. بنابراین، ما با برداری کردن مناسب کلمات به قدرت مناسب پردازش هم دست پیدا می کنیم و این بخش اهمیت بالایی دارد. لایه Embedding اصولاً یک ماتریس است که دو بعد دارد و تعداد سطرهای آن معادل تعداد کلمات



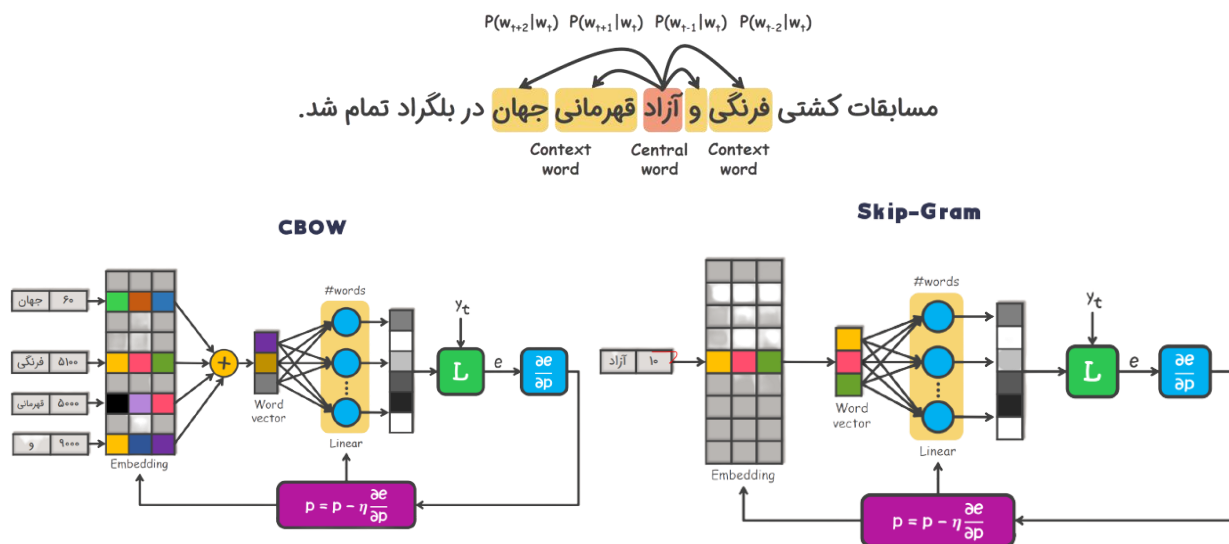
شکل ۲: فرآیند تبدیل کلمه به بردار.

دیکشنری است و هر سطر معادل یک کلمه است. و هم چنین تعداد ستون ها نشان دهنده طول برداری است که می خواهیم برای هر کلمه داشته باشیم. مثالی از این موضوع در شکل ۳ نشان داده شده است. این لایه شبیه به همان لایه خطی است و این ماتریس هم قابل آموزش است و باید آموزش داده شود تا مقادیر آن معنا دارد شود (مثلاً کلمات مرد و زن به هم شبیه شوند).



شکل ۳: لایه Embedding.

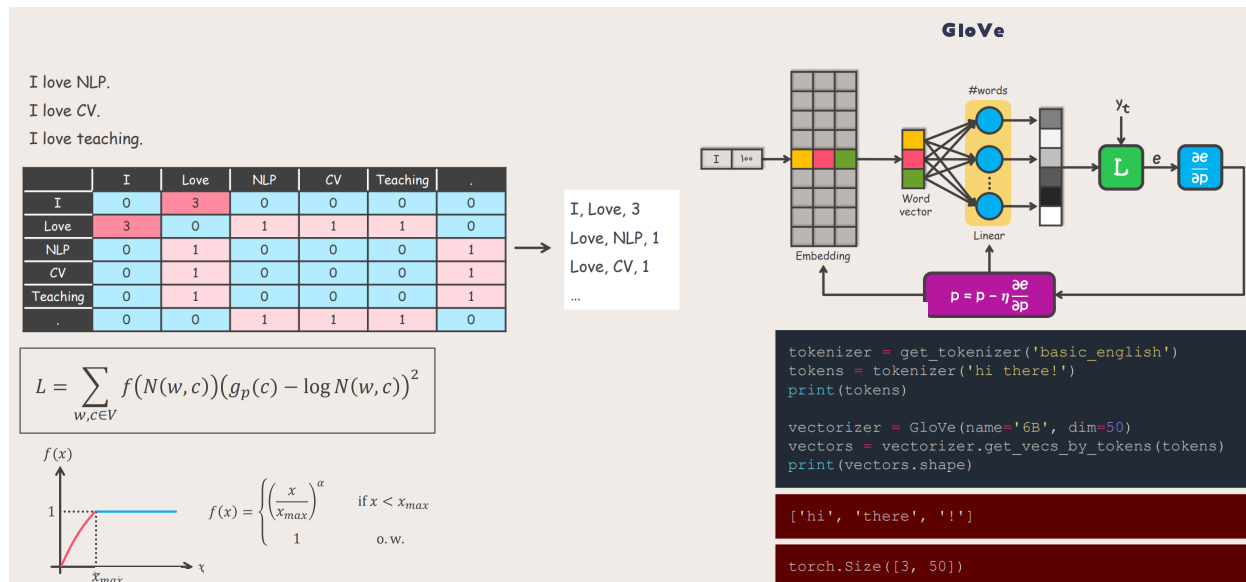
در ادامه بد نیست به نحوه فهم کلمات جدید برای خودمان به عنوان یک انسان توجه کنیم. ما وقتی معنای کلمه‌ای را نمی‌دانیم و آن را در جملاتی مشاهده می‌کنیم، بر اساس کلمات قبل و بعدی که از آن می‌بینیم در مورد مفهوم آن کلمه ناشناخته فهمی پیدا می‌کنیم. بر اساس همین ایده یک الگوریتم قدیمی به نام Word2Vec ایجاد شده که در آن می‌توان مفهوم یک کلمه مرکزی (Central Word) را بر اساس پنجره‌ای از کلمات همسایه (Context Word) را حدس زد. این کار با الگوریتم CBOW انجام می‌شود. این رابطه دوطرفه است و از کلمه مرکزی هم می‌توان به کلمات همسایه رسید. این کار با الگوریتم Skip-Gram انجام می‌شود. شمایی از این الگوریتم در شکل ۴ نشان داده شده است. در این الگوریتم کلمات کم‌کاربرد مانند حروف ربط را معمولاً حذف می‌کنیم. می‌توان از اوزان آماده ماتریس نهفته هم استفاده کرد. Glove هم یکی از روش‌هایی است که می‌توانیم با آن



شکل ۴: Word2Vec.

یک لایه Embedding معنادار و خوبی داشته باشیم. اوزان پیش‌آموزش دیده آن هم در پای‌تورج آماده استفاده است. تفاوت کوچک این روش با Word2Vec این است که در این جا زوج کلماتی که می‌خواهیم بسازیم بر اساس یک ماتریس هم‌رخداد یا Co-occurrence Matrix ساخته می‌شود. فرض کنید یکسری جمله داشته باشیم و در آن می‌خواهیم بررسی کنیم کلمات چندبار

و به چه صورتی در کنار هم قرار گرفته‌اند. نمونه‌ای از این فرآیند ساخت ماتریس هم‌رخداد در **شکل ۵** نشان داده شده است. این ماتریس اطلاعات ارزشمندی به فرآیند آموزش اضافه می‌کند. ساختار آموزش این روش کاملاً مشابه قبل است اما تابع اتلاف آن به جای آنتروپی متقابل بر پایه MSE (معادله داخل **شکل ۵**) است. در این معادله زوج کلمات مدنظر ما هستند. N نمایانگر همان ماتریس هم‌رخداد است و g هم خروجی کلاسیفایر پیش‌بینی و لایه نهفته است. در نهایت یک وزن هم در این رابطه داریم. این وزن عمده تمرکزش بر واژگانی است که ممکن است اتفاقی چندبار کنار هم قرار گرفته باشند و اهمیت بالایی هم نداشته باشند. درواقع این وزن می‌گوید این زوج کلمات باید مثلاً حداقل ۱۰۰ بار کنار هم دیده شوند که اهمیت بالایی پیدا کنند.



شکل ۵: GloVe.

نقطه ضعف پررنگ روش‌هایی که تاکنون اشاره کردیم فهم معنای کلمات چندمعنا است و لایه Embedding به‌ازای آن کلمه بدون توجه به معنای آن تنها یک بردار واحد نسبت می‌دهد. روش شبکه ELMO بر پایه LSTM سعی می‌کند برای هر کلمه با توجه به جمله‌ای که در آن قرار گرفته یک بردار معنادار بسازد. در واقع ما به این شبکه جمله می‌دهیم و می‌گوییم به دنبال معنا و بردار معنادار کلمات هستیم. اما بعد از آمدن ترنسفورمرها شبکه‌ای بر پایه ترنسفورمر به نام BERT آمد که سعی داشت تمام مشکلات مانند همین مشکل چندمعنایی را حل کند (BERT Word Embeddings).

۲.۱ پاسخ قسمت ۲

مقاله در مورد مجموعه داده این‌گونه بیان می‌کند که به منظور آموزش و آزمایش مدل‌های مختلف، از یک مجموعه داده ناهمگن استفاده شده است. این مجموعه داده شامل سوالاتی است که به صورت دستی ساخته شده و همچنین سوالاتی که توسط TREC و USC منتشر شده‌اند. این مجموعه داده شامل ۵۵۰۰ سوال در زبان انگلیسی در مجموعه آموزش و ۵۰۰ سوال در مجموعه آزمایش است. هر سوال دارای یک برچسب است که پاسخ را در یک سلسله مراتب دو سطحی دسته‌بندی می‌کند. سطح بالاتر شامل شش کلاس اصلی است (کوتاه‌نویسی، موجودیت، توصیف، انسان، مکان، عدد)، هر یک از این کلاس‌ها در نتیجه به زیرکلاس‌های متمایزی تقسیم می‌شوند که با هم در سطح دوم سلسله مراتب قرار می‌گیرند (مانند: انسان/گروه، انسان/شخصیت، مکان/شهر، مکان/کشور و غیره). به طور کلی، از ترکیب همه دسته‌ها و زیردسته‌ها، ما ۵۰ برچسب مختلف برای دسته‌بندی پاسخ

به سوالات در اختیار داریم. یک مثال که از مجموعه داده استخراج شده است به صورت زیر است (برچسب با حروف توپر و سوال به صورت ایتالیک درج شده است):

HUMAN:ind *Who was the 16th President of the United States?*

همانطور که مشاهده می شود، برچسب از دو بخش تشکیل شده است که توسط علامت ":" از هم جدا شده اند و به ترتیب کلاس اصلی و زیرکلاس را نمایش می دهند. در این مثال، کلاس اصلی نشان می دهد که پاسخ باید یک شخص یا یک گروه را مشخص کند، در حالی که زیرکلاس نشان می دهد که ما می خواهیم یک فرد خاص را (به طور طبیعی از طریق نام او) شناسایی کنیم. لازم به ذکر است که بازنمایی ارائه شده توسط GloVe شامل تمامی کلمات (۹۱۲۳ کلمه) موجود در مجموعه داده می شود و بدین ترتیب نیازی به کارهای بیشتر در این زمینه نیست. با این حال، پس از استخراج داده ها از مجموعه داده، ابتدا یک پیش پردازش اولیه انجام می شود که شامل پاکسازی داده ها از هر نویسه و علامت نگارشی غیر مفید (به استثنای علامت سوال) می شود. در جلوتر مشاهده می شود که برای عمل به این گفته از دستور زیر استفاده شده است تا علامت سوال باقی بماند:

```
text = re.sub(r'[^\a-zA-Z0-9\?]+', ' ', text)
```

علاوه بر این، تمام کلمات موجود در رشته ها به حروف کوچک تبدیل می شوند تا از تکرار کدگذاری های چندگانه برای کلمات و عبارات یکسان جلوگیری شود. حال در گام اول و برای راحتی، مجموعه داده را در [گوگل درایو](#) بارگذاری می کنیم. هم چنین، فراخوانی این فایل بدون نیاز به Mount کردن و استفاده از gdown در محیط گوگل کولب ممکن است با خطا مواجه شود که با استفاده از دستورات زیر مشکل حل می شود (منبع). در ادامه هم دستوراتی برای دریافت اطلاعات مربوط به نهفته سازی Glove نوشته ایم.

```
1 !pip install --upgrade --no-cache-dir gdown
2 !gdown 1lDpsLB-erPd4rvvRguui6i06h1hTKE1K
3 !gdown 1HctYMsZ-V7t7ipdL0005S-WNWH5X_14b
4 !gdown 1nZ-JnAIQ5fku0FhIFUHCfKGuCtEyuNNB
5
6 # Load pre-trained GloVe embeddings
7 !wget http://nlp.stanford.edu/data/glove.6B.zip
8 !unzip glove*.zip
9
10 nltk.download('punkt')
```

در ادامه دستورات مختلفی را برای آماده سازی و پیش پردازش داده ها می نویسیم. در این دستورات، متن داده های آموزشی و آزمون پاکسازی می شود. ابتدا تمامی کاراکترهای خاص و ناخواسته از متن حذف می شوند و سپس متن به حروف کوچک تبدیل می شود. در ادامه، متن پاکسازی شده توسط تابع word_tokenize توکن بندی می شود. این فرایند شامل تقسیم متن به واحدهای جداگانه که معمولاً کلمات هستند می باشد. سپس، باید مجموعه ی کلمات متوقف یا متوقف را دریافت کنیم. کلمات متوقف یا stop words در پردازش زبان طبیعی به کلماتی اشاره دارند که اغلب در متن ها به صورت تکراری و بدون ارزش اطلاعاتی ظاهر می شوند. این کلمات عموماً از نظر معنایی کم اهمیت هستند و بیشتر برای ایجاد جریان و قواعد گرامری در جملات استفاده می شوند. مثال هایی از کلمات متوقف شامل "the"، "is"، "and"، "a"، "an" و "in" می باشند. در پردازش متن، حذف کلمات متوقف می تواند منجر به بهبود عملکرد الگوریتم ها و مدل ها در وظایفی مانند تحلیل احساسات، تشخیص موضوع و دسته بندی متن شود. با حذف کلمات متوقف، تمرکز بیشتری بر کلمات کلیدی و اطلاعات معنایی موجود در متن داریم. با استفاده از nltk.download('stopwords')، منابع مربوط به کلمات متوقف دانلود می شوند. سپس کلمات متوقف انگلیسی در stop_words قرار می گیرند. در نهایت، این کلمات از متن پاکسازی شده حذف می شوند و نتیجه در ستون filtered_text قرار می گیرد. در ادامه، متن فیلتر شده و برچسب های داده های آموزشی و آزمون را از داده های اصلی استخراج می کنیم و آن ها را به

ترتیب در متغیرهای مربوط به خود قرار می‌دهیم. ابتدا تابع Tokenizer را با تعداد کلمات مجاز برابر با ۴۰۰۰۰۰ ایجاد می‌کنیم و با استفاده از fit_on_texts متن آموزشی را برای توکن‌بندی آماده می‌کنیم. سپس با استفاده از texts_to_sequences، متن‌های آموزشی و آزمون را به دنباله‌های عددی تبدیل می‌کنیم و در متغیرهای X_train و X_test ذخیره می‌کنیم. حداکثر طول دنباله را برابر با بیشترین طول دنباله در X_train در نظر می‌گیریم و سپس با استفاده از pad_sequences، دنباله‌های عددی را پُر می‌کنیم تا طول یکسانی داشته باشند و در متغیرهای مربوطه ذخیره می‌کنیم. تعداد دسته‌ها را در متغیر num_classes ذخیره می‌کنیم. سپس با استفاده از np.eye، برچسب‌های آموزشی و آزمون را به بردارهای دودویی تبدیل می‌کنیم و در متغیرهای y_train و y_test ذخیره می‌کنیم. در ادامه کار، مسیر فایل glove_path را با مسیر فایل Embeddings GloVe مدنظر خود جایگزین می‌کنیم. سپس فایل را خوانده و در آرایه embedding_matrix ذخیره می‌کنیم. در این آرایه، بردارهای embedding برای کلمات موجود در word_index قرار می‌گیرند. در صورتی که برداری برای یک کلمه موجود نباشد، مقدار صفر به آن اختصاص داده می‌شود. لازم به ذکر است که برای نرمال‌سازی از تابع normalize_text استفاده می‌شود. دستورات به شرح زیر هستند:

```
1 ## Method1
2
3 import pandas as pd
4 import numpy as np
5 import nltk
6 from nltk.corpus import stopwords
7 from nltk.tokenize import word_tokenize
8 from sklearn.model_selection import train_test_split
9 from keras.preprocessing.text import Tokenizer
10 from tensorflow.keras.preprocessing.sequence import pad_sequences
11 import warnings
12 # Ignore all warnings
13 warnings.filterwarnings("ignore")
14
15 train_data = pd.read_csv('train.csv')
16 test_data = pd.read_csv('test.csv')
17
18
19 # Clean the text by removing special characters and converting to lowercase
20 train_data['cleaned_text'] = train_data['text'].str.replace('[^\w\s?]', '').str.lower()
21 test_data['cleaned_text'] = test_data['text'].str.replace('[^\w\s?]', '').str.lower()
22
23 # Tokenize the text
24 train_data['tokenized_text'] = train_data['cleaned_text'].apply(word_tokenize)
25 test_data['tokenized_text'] = test_data['cleaned_text'].apply(word_tokenize)
26
27 # Remove stopwords
28 nltk.download('stopwords')
29 stop_words = set(stopwords.words('english'))
30 train_data['filtered_text'] = train_data['tokenized_text'].apply(lambda tokens: [word for word in
    tokens if word not in stop_words])
31 test_data['filtered_text'] = test_data['tokenized_text'].apply(lambda tokens: [word for word in
```

```

tokens if word not in stop_words])
32
33
34 # Get the filtered text and labels
35 X_train = train_data['filtered_text'].values
36 y_train = train_data['label-fine'].values
37 X_test = test_data['filtered_text'].values
38 y_test = test_data['label-fine'].values
39
40 # Tokenize the text and convert to sequences
41 tokenizer = Tokenizer(num_words=400000)
42 tokenizer.fit_on_texts(X_train)
43
44 X_train = tokenizer.texts_to_sequences(X_train)
45 X_test = tokenizer.texts_to_sequences(X_test)
46
47 # Pad the sequences to have the same length
48 max_sequence_length = max(len(seq) for seq in X_train)
49 X_train = pad_sequences(X_train, maxlen=max_sequence_length)
50 X_test = pad_sequences(X_test, maxlen=max_sequence_length)
51
52 # Convert labels to categorical
53 num_classes = train_data['label-fine'].nunique()
54 y_train = np.eye(num_classes)[y_train]
55 y_test = np.eye(num_classes)[y_test]
56
57 # Define function for text normalization
58 def normalize_text(text):
59     text = re.sub(r'[^a-zA-Z0-9\?]+', ' ', text)
60     text = text.lower()
61     return text
62
63 # Normalize the filtered text
64 X_train = [normalize_text(text) for text in X_train]
65 X_test = [normalize_text(text) for text in X_test]
66
67 # Replace 'glove_path' with the path to your GloVe embeddings file
68 glove_path = '/content/glove.6B.300d.txt'
69 embeddings_index = {}
70 with open(glove_path, encoding='utf8') as f:
71     for line in f:
72         values = line.split()
73         word = values[0]
74         coefs = np.asarray(values[1:], dtype='float32')
75         embeddings_index[word] = coefs

```

```

76
77 # Create an embedding matrix
78 embedding_dim = 300
79 word_index = tokenizer.word_index
80 num_words = min(400000, len(word_index) + 1)
81 embedding_matrix = np.zeros((num_words, embedding_dim))
82 for word, i in word_index.items():
83     if i >= num_words:
84         continue
85     embedding_vector = embeddings_index.get(word)
86     if embedding_vector is not None:
87         embedding_matrix[i] = embedding_vector
88
89
90 ## Method2
91
92 # Import necessary libraries
93 import pandas as pd
94 import numpy as np
95 import re
96 import nltk
97 from nltk.tokenize import word_tokenize
98 from tensorflow.keras.preprocessing.text import Tokenizer
99 from tensorflow.keras.preprocessing.sequence import pad_sequences
100 from sklearn.preprocessing import LabelEncoder
101
102 # Load dataset
103 train_data = pd.read_csv('train.csv')
104 test_data = pd.read_csv('test.csv')
105
106 # Read 'QA_data.csv' with 'errors' parameter and pass the file object to 'pd.read_csv()'
107 with open('QA_data.csv', 'r', encoding='utf-8', errors='replace') as file:
108     qa_data = pd.read_csv(file)
109
110 # Define function for text normalization
111 def normalize_text(text):
112     text = re.sub(r'[^a-zA-Z0-9\?]+', ' ', text)
113     text = text.lower()
114     return text
115
116 # Normalize text
117 train_data['text'] = train_data['text'].apply(normalize_text)
118 test_data['text'] = test_data['text'].apply(normalize_text)
119 qa_data['text'] = qa_data['text'].apply(normalize_text)
120

```

```

121 # Tokenize the text
122 nltk.download('punkt')
123 train_data['tokens'] = train_data['text'].apply(word_tokenize)
124 test_data['tokens'] = test_data['text'].apply(word_tokenize)
125 qa_data['tokens'] = qa_data['text'].apply(word_tokenize)
126
127 # Load pre-trained GloVe embeddings
128 !wget http://nlp.stanford.edu/data/glove.6B.zip
129 !unzip glove*.zip
130
131 # Load GloVe embeddings into a dictionary
132 embeddings_index = {}
133 with open('glove.6B.300d.txt') as f:
134     for line in f:
135         values = line.split()
136         word = values[0]
137         coefs = np.asarray(values[1:], dtype='float32')
138         embeddings_index[word] = coefs
139
140 # Prepare tokenizer
141 tokenizer = Tokenizer()
142 tokenizer.fit_on_texts(train_data['tokens'])
143
144 # Convert tokens to sequences
145 train_data['sequences'] = tokenizer.texts_to_sequences(train_data['tokens'])
146 test_data['sequences'] = tokenizer.texts_to_sequences(test_data['tokens'])
147 qa_data['sequences'] = tokenizer.texts_to_sequences(qa_data['tokens'])
148
149 # Pad sequences
150 maxlen = max(train_data['sequences'].apply(len))
151 train_padded_sequences = pad_sequences(train_data['sequences'], maxlen=maxlen)
152 test_padded_sequences = pad_sequences(test_data['sequences'], maxlen=maxlen)
153 qa_padded_sequences = pad_sequences(qa_data['sequences'], maxlen=maxlen)
154
155 # Encode labels
156 encoder = LabelEncoder()
157 encoder.fit(train_data['label-coarse'])
158 train_data['encoded_labels'] = encoder.transform(train_data['label-coarse'])
159 test_data['encoded_labels'] = encoder.transform(test_data['label-coarse'])
160 qa_data['encoded_labels'] = encoder.transform(qa_data['label-coarse'])
161
162 # Create an embedding matrix
163 embedding_dim = 300
164 vocab_size = len(tokenizer.word_index) + 1
165 embedding_matrix = np.zeros((vocab_size, embedding_dim))

```

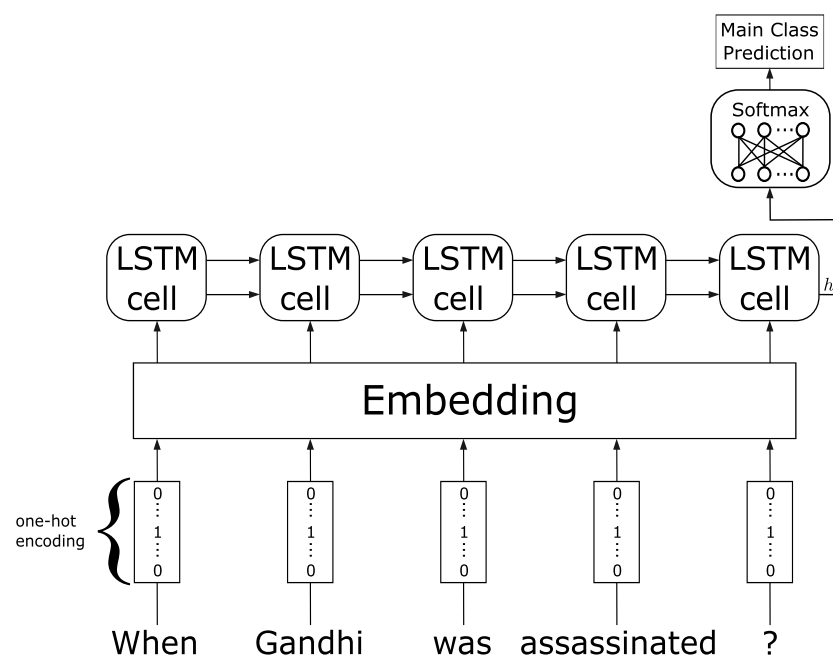
```

166 for word, i in tokenizer.word_index.items():
167     embedding_vector = embeddings_index.get(word)
168     if embedding_vector is not None:
169         embedding_matrix[i] = embedding_vector

```

۳.۱ پاسخ قسمت ۳

مقاله در مورد مدل اول این گونه بیان می کند که تمام دنباله ها را که به شبکه وارد می شوند، به یک بردار با ابعاد ثابت، که ابعاد آن برابر با تعداد دسته بندی های اصلی است نگاشت می کنیم. به عبارت دیگر در این روش، هر کلمه تشکیل دهنده سوال (طبق ترتیب زمانی و پس از انجام نگاشت از طریق لایه تعبیه) وارد شبکه LSTM شود و حالت و طبقه اصلی نهایی آن از طریق یک شبکه عصبی کلاسیک با فعال سازی Softmax در یکی از شش کلاس پیش بینی اصلی پیش بینی می شود. الگوریتم از Backpropagation through time (BPTT) و تابع هزینه Categorical Cross-Entropy استفاده می کند. نمایی از این مدل در شکل ۶ نشان داده شده است. برای پیاده سازی دستورات مربوط به این مدل، ابتدا داده های آموزش و آزمون که از مرحله قبل ایجاد شده اند را



شکل ۶: نمایی از مدل اول.

به صورت کامل معین و آماده می کنیم. تابع encode_labels در زبان برنامه نویسی پایتون، وظیفه تبدیل برچسب ها به بردارهای One-Hot را دارد. این تابع دو ورودی دریافت می کند: labels که برچسب های ورودی را نشان می دهد و num_classes که تعداد کلاس ها را نمایش می دهد. در داخل تابع، با استفاده از تابع tf.keras.utils.to_categorical، برچسب ها به بردارهای One-Hot تبدیل می شوند (مطابق خواست مقاله در شکل ۶). هر برچسب به یک بردار با طول num_classes تبدیل می شود. در این بردار، مکان متناظر با شماره کلاس برچسب، با مقدار ۱ و سایر مکان ها با مقدار ۰ پر می شود. به عبارت دیگر، هر برچسب به صورت یک بردار با مقدار ۱ در مکان متناظر با کلاس و مقادیر ۰ در سایر مکان ها نمایش داده می شود.

برای تعریف مدل یک تابع با نام `create_lstm_model` تعریف می‌کنیم که یک مدل LSTM را ایجاد می‌کند. مدل از ساختار Sequential استفاده می‌کند، که به ترتیب لایه‌های Embedding، LSTM و Dense را اضافه می‌کند. لایه Embedding از واژگان و جملات ورودی استفاده می‌کند و وزن‌های جدول تعبیه (`embedding_matrix`) را استفاده می‌کند. وزن‌های جدول تعبیه از قبل آموزش دیده‌اند و برای نمایش برداری واژگان و تبدیل واژگان به بردارهای عددی استفاده می‌شوند. ورودی لایه Embedding دارای اندازه واژگان (`vocab_size`) و بعد بردارهای تعبیه (`embedding_dim`) است و همچنین طول ورودی (`input_length`) را نیز تعیین می‌کنیم. لایه بعدی یک LSTM است که تعداد واحدهای LSTM را با پارامتر `h_dim` تعیین می‌کند. این لایه به صورت توالی بازگشتی (`return_sequences=True`) است تا مطابق شکل ۶ خروجی‌ها را به لایه بعدی انتقال دهد. تابع فعال‌سازی لایه LSTM هم تانژانت هیپربولیک (`tanh`) در نظر گرفته‌ایم. سپس، می‌توان یک لایه دیگر LSTM با تعداد واحدهای LSTM مشابه قبلی اضافه کرد. این لایه در حالت پیش‌فرض تک لایه LSTM است که خروجی نهایی را تولید می‌کند. در نهایت، یک لایه Dense با تابع فعال‌سازی softmax برای تولید احتمالات کلاس‌ها اضافه می‌شود. تعداد کلاس‌ها با پارامتر `num_classes` مشخص می‌شود. مدل با استفاده از تابع هزینه `categorical_crossentropy` و بهینه‌ساز Adam (با نرخ یادگیری ۰.۰۰۱) تنظیم می‌شود. معیارهای دقت و اتلاف نیز برای مدل مشخص شده است. در نهایت، مدل ساخته شده با استفاده از توابع `compile` و `return` آماده استفاده و آموزش می‌شود. سپس فرآیندها را تعیین می‌کنیم.

در ادامه توابعی برای رسم نمودارها و محاسبات لازم برای حل سوال ایجاد می‌کنیم. تابع `save_plots_to_pdf` را برای ذخیره نمودارها به صورت فایل PDF نوشته‌ایم. در این تابع نمودارهای دقت و اتلاف برای داده‌های آموزش و اعتبارسنجی رسم شده و پارامترهای مهم مدنظر سوال (مانند `h_dim`) در قسمت عنوان نوشته می‌شوند. نمودارهای هم ذخیره می‌شوند و هم در محیط کولب نمایش داده می‌شوند. با استفاده از `plt.close`، منابع گرافیکی آزاد می‌شوند و نمودارها بسته می‌شوند. تابع `save_confusion_matrix_to_pdf` هم برای رسم ماتریس درهم‌ریختگی به شکلی مناسب و به صورتی که در آن پارامترهای مهم مدنظر سوال (مانند `h_dim`) نشان داده شود نوشته شده است. با نوشتن این توابع به سراغ نوشتن بخش اصلی کد می‌رویم. این قسمت با استفاده از مجموعه‌ای از پارامترهای انتخاب شده و با استفاده از مدل LSTM پیش‌تر تعریف شده، مدل‌های مختلف را آموزش می‌دهد و نتایج عملکرد آن‌ها را گزارش می‌دهد. در بخش آموزش این کد، برای هر مقدار ابعاد `h`، یک مدل LSTM با استفاده از تابع `create_lstm_model` ایجاد می‌شود و با استفاده از داده‌های آموزش و اعتبارسنجی، مدل آموزش داده می‌شود. در بخش ارزیابی هم دقت مدل روی مجموعه‌های آموزش و آزمون محاسبه می‌شود. همچنین، از معیارهای دیگر اشاره شده در صورت سوال هم برای ارزیابی عملکرد مدل استفاده می‌شود. در نهایت نمودارهای مختلف از جمله نمودارهای دقت، ماتریس‌های درهم‌ریختگی و معیارهای دیگر برای هر مدل در یک فایل PDF ذخیره می‌شوند. قبل از ادامه به ابعاد بعدی `h`، برنامه طوری تنظیم شده که ۲ ثانیه منتظر می‌ماند. این کد به صورت سیستماتیک و در چندین مرحله، مدل LSTM پایه را با ابعاد مختلف آموزش می‌دهد و عملکرد آن را بررسی می‌کند. همچنین، نمودارها و ماتریس‌های درهم‌ریختگی برای هر مدل ذخیره می‌کند تا بتوان ارزیابی دقیق‌تری انجام داد. با ارائه این توضیحات تمام دستورات مربوطه در برنامه ۱ آورده شده است. به دلیل عدم شفافیت کامل دستورات مورد استفاده مقاله از قبیل عدم مشخص بودن اندازه دسته، تعداد دوره، در نظر گرفتن داده اعتبارسنجی و غیره، از تجربیات اندوخته تمرین‌های گذشته برای رسیدن به یک نتیجه مناسب و قابل قبول استفاده کرده‌ایم. مثلاً مدل را در چند حالت مختلف تشکیل داده‌ایم و بهترین نتیجه را این‌جا ذکر کرده‌ایم. هم‌چنین دو حالت در نظر گرفته‌ایم. یک حالت نزدیک به آن چیزی که احتمالاً مقاله پیش رفته (بدون داده اعتبارسنجی و با حدود ۵۰ دوره آموزش) و یکی هم با داده اعتبارسنجی. دستورات آورده شده در برنامه ۱ حالت کامل‌تر و با داده اعتبارسنجی است و برای جلوگیری از ازدحام بیش‌تر دستورات حالتی که در آن داده اعتبارسنجی در نظر نگرفته‌ایم آورده نشده است. برای قسمت لایه Embedding آموزش را غیرفعال می‌کنیم چون ممکن است آموزش و ایجاد ارتباط بین کلمات نسبت به حالت قدرتمند پیش آموزش دیده موجود خوب انجام نپذیرد.

Program 1: Model One, Training, and Evaluation

```

1 import tensorflow as tf
2 from tensorflow.keras.models import Sequential
3 from tensorflow.keras.layers import Embedding, LSTM, Dense
4 from tensorflow.keras.optimizers import Adam
5 import numpy as np
6 import pandas as pd
7 import matplotlib.pyplot as plt
8 from sklearn.metrics import f1_score, precision_score, recall_score, confusion_matrix
9 import seaborn as sns
10 from matplotlib.backends.backend_pdf import PdfPages
11 import time
12
13 # Set random seeds for reproducibility
14 np.random.seed(42)
15 tf.random.set_seed(42)
16
17 # Prepare the data
18 X_train = train_padded_sequences
19 y_train = train_data['encoded_labels']
20 X_test = test_padded_sequences
21 y_test = test_data['encoded_labels']
22
23 # Function to encode labels as one-hot vectors
24 def encode_labels(labels, num_classes):
25     return tf.keras.utils.to_categorical(labels, num_classes=num_classes)
26
27 from keras.regularizers import l2
28
29 # Function to create the LSTM model
30 def create_lstm_model(input_length, h_dim, num_classes):
31     model = Sequential()
32     model.add(Embedding(input_dim=vocab_size, output_dim=embedding_dim, input_length=input_length
33         , weights=[embedding_matrix], trainable=False))
34     model.add(LSTM(h_dim, return_sequences=True, activation='tanh'))
35     model.add(LSTM(h_dim, activation='tanh'))
36     model.add(Dense(num_classes, activation='softmax', kernel_regularizer=l2(0.001)))
37     model.compile(loss='categorical_crossentropy', optimizer=Adam(learning_rate=0.001), metrics=[
38         'accuracy'])
39     return model
40
41 # def create_lstm_model(input_length, h_dim, num_classes):
42 #     model = Sequential()
43 #     model.add(Embedding(input_dim=vocab_size, output_dim=embedding_dim, input_length=
44 #         input_length, weights=[embedding_matrix], trainable=False))
45 #     model.add(LSTM(h_dim, activation='tanh'))

```



```

42 #     model.add(Dense(num_classes, activation='softmax'))
43 #     model.compile(loss='categorical_crossentropy', optimizer=Adam(learning_rate=0.001), metrics
    =['accuracy'])
44 #     return model
45
46 # Function to save plots to PDF
47 def save_plots_to_pdf(filename):
48     with PdfPages(filename) as pdf:
49         plt.figure(figsize=(10, 4))
50
51         plt.subplot(1, 2, 1)
52         plt.plot(history.history['accuracy'], label='Training')
53         plt.xlabel('Epoch')
54         plt.ylabel('Accuracy')
55         plt.legend()
56         plt.title(f'Accuracy (h_dim = {h_dim})')
57
58         plt.subplot(1, 2, 2)
59         plt.plot(history.history['loss'], label='Training')
60         plt.xlabel('Epoch')
61         plt.ylabel('Loss')
62         plt.legend()
63         plt.title(f'Loss (h_dim = {h_dim})')
64
65         pdf.savefig(bbox_inches='tight')
66
67         # Show the plots in Colab output
68         plt.show()
69
70     plt.close()
71
72 # Function to save confusion matrix to PDF
73 def save_confusion_matrix_to_pdf(y_true, y_pred, filename):
74     cm = confusion_matrix(y_true, y_pred)
75     plt.figure(figsize=(6, 6))
76     sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', cbar=False)
77     plt.xlabel('Predicted label')
78     plt.ylabel('True label')
79     plt.title(f'Confusion matrix (h_dim = {h_dim})')
80     plt.savefig(filename, bbox_inches='tight')
81
82     # Show the plot in Colab output
83     plt.show()
84     plt.close()
85

```

```

86 # Set hyperparameters
87 h_dimensions = [25, 50, 75, 100]
88 epochs = 50
89 num_classes = 6
90 batch_size = 64
91
92 # Iterate over h_dimensions and train the models
93 for h_dim in h_dimensions:
94     print(f"Training model with h_dim = {h_dim}")
95     model = create_lstm_model(maxlen, h_dim, num_classes)
96
97     # Encode labels as one-hot vectors
98     y_train_encoded = encode_labels(y_train, num_classes)
99
100     history = model.fit(X_train, y_train_encoded, epochs=epochs, batch_size=batch_size, verbose
    =2)
101     train_accuracy = model.evaluate(X_train, encode_labels(y_train, num_classes), verbose=0)[1] *
    100
102     test_accuracy = model.evaluate(X_test, encode_labels(y_test, num_classes), verbose=0)[1] *
    100
103     print(f"Training set accuracy: {train_accuracy:.2f}%")
104     print(f"Test set accuracy: {test_accuracy:.2f}%")
105
106     y_pred_train = np.argmax(model.predict(X_train), axis=1)
107     y_pred_test = np.argmax(model.predict(X_test), axis=1)
108
109     train_f1 = f1_score(y_train, y_pred_train, average='weighted')
110     test_f1 = f1_score(y_test, y_pred_test, average='weighted')
111
112     train_precision = precision_score(y_train, y_pred_train, average='weighted')
113     test_precision = precision_score(y_test, y_pred_test, average='weighted')
114
115     train_recall = recall_score(y_train, y_pred_train, average='weighted')
116     test_recall = recall_score(y_test, y_pred_test, average='weighted')
117
118     print(f"Training set F1-score: {train_f1:.2f}")
119     print(f"Test set F1-score: {test_f1:.2f}")
120
121     print(f"Training set Precision: {train_precision:.2f}")
122     print(f"Test set Precision: {test_precision:.2f}")
123
124     print(f"Training set Recall: {train_recall:.2f}")
125     print(f"Test set Recall: {test_recall:.2f}")
126
127     save_plots_to_pdf(f"model_{h_dim}_plots.pdf")

```

```

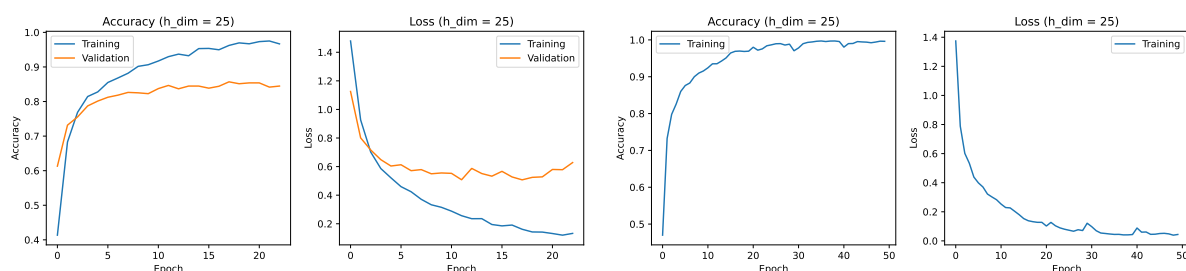
128
129 print("Confusion matrix (training set):")
130 save_confusion_matrix_to_pdf(y_train, y_pred_train, f"model{h_dim}confusion_matrix_train.pdf"
131 )
132
133 print("Confusion matrix (test set):")
134 save_confusion_matrix_to_pdf(y_test, y_pred_test, f"model{h_dim}confusion_matrix_test.pdf")
135
136 print("\n")
137
138 time.sleep(2) # Wait for 2 seconds before proceeding to the next h_dim

```

نتایج در حالت نزدیک به مقاله، به صورتی است که در جدول ۱ و شکل ۷ تا شکل ۱۴ نشان داده شده است. همان طور که مشاهده می شود. تفکیک ۶ کلاس اصلی با دقت قابل قبولی انجام شده است. هرچند برخی نتایج نشان داده شده برای داده های آموزش اهمیت و معنای خاصی ندارند (مانند ماتریس درهم ریختگی و...) آن ها را هم نشان دادیم تا صرفاً به برخی نکات مانند چالش در تشخیص برچسب ۲ به دلیل کم تر بودن داده در آن دسته اشاره کنیم.

جدول ۱: نتایج مربوط به مدل اول

Test Set (%)				Training Set (%)				h Dimension
F1-score	Recall	Precision	Accuracy	F1-score	Recall	Precision	Accuracy	
۸۹	۸۹	۸۹	۸۹.۰۰	۱۰۰	۱۰۰	۱۰۰	۹۹.۸۳	25
۹۰	۹۰	۹۱	۹۰.۲۰	۱۰۰	۱۰۰	۱۰۰	۹۹.۹۶	50
۹۱	۹۲	۹۲	۹۱.۶۰	۱۰۰	۱۰۰	۱۰۰	۹۹.۹۸	75
۹۱	۹۱	۹۱	۹۱	۱۰۰	۱۰۰	۱۰۰	۹۹.۹۸	100

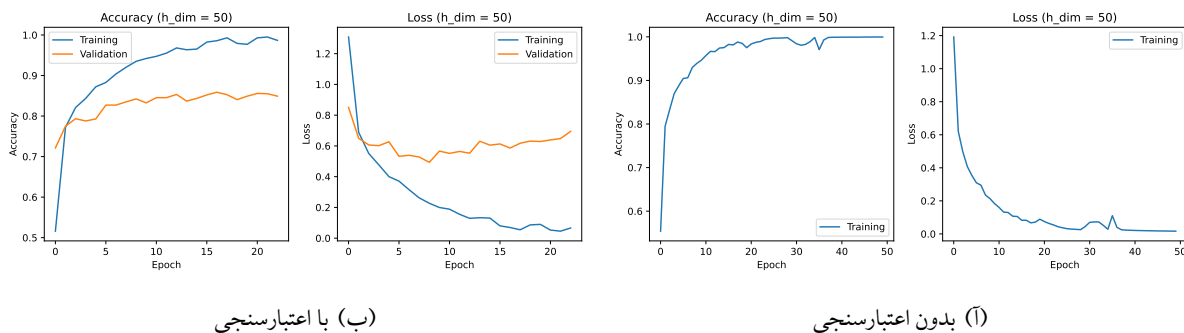


(ب) با اعتبارسنجی

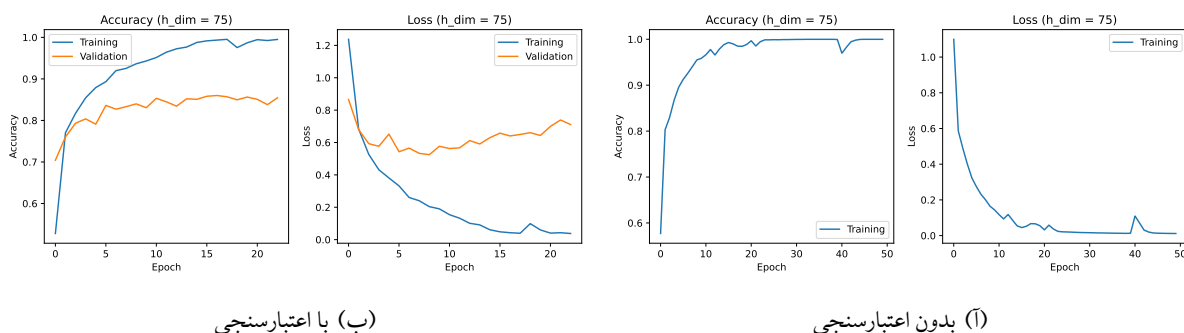
(آ) بدون اعتبارسنجی

شکل ۷: نمودار دقت و اتلاف - مدل اول - h=25.

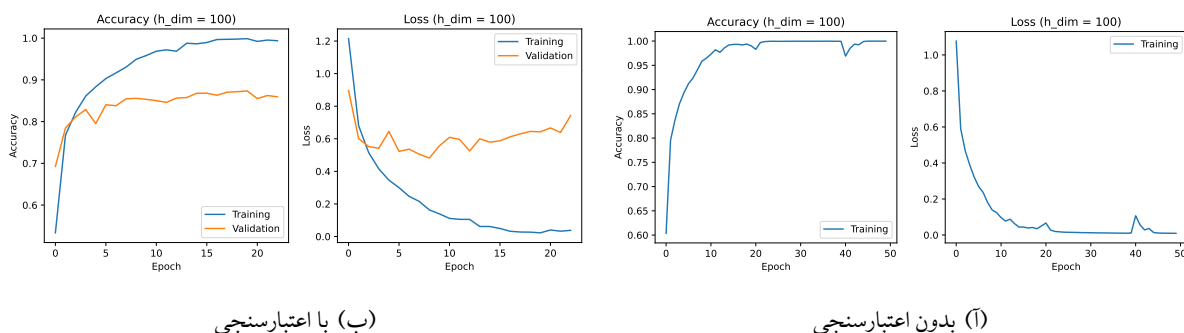
مقاله این گونه بیان می کند که نتایج خوبی که در پیش بینی کلاس اصلی به دست آمد، انگیزه ای برای توسعه مدل با حفظ بخش اول و اصلی آن شد. در بخش دوم، پیش بینی زیرکلاس نیز در نظر گرفته شده است. زیرکلاس باید به عنوان یک ویژگی تخصصی از کلاس اصلی تأثیر بگیرد. ایده اصلی مقاله برای بخش دوم این است که یک عنصر دیگر به انتهای سوال اضافه شود. این عنصر پرکننده اطلاعات مفیدی ندارد ولی تنها برای ارزیابی خروجی پس از خروجی متناظر با آخرین عنصر سوال لازم است. برای این کار، خروجی قبل از آخرین عنصر سوال مانند حالت قبلی می تواند در نظر گرفته شود، در حالی که خروجی آخر تنها



شکل ۸: نمودار دقت و اتلاف - مدل اول - $h=50$.

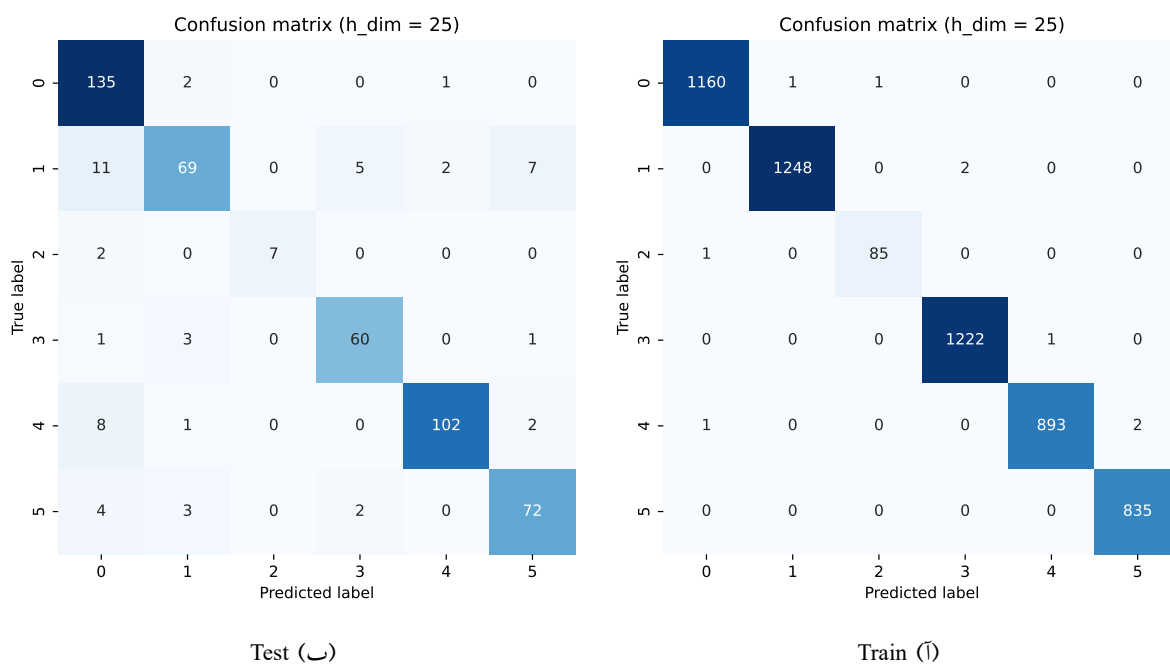


شکل ۹: نمودار دقت و اتلاف - مدل اول - $h=75$.

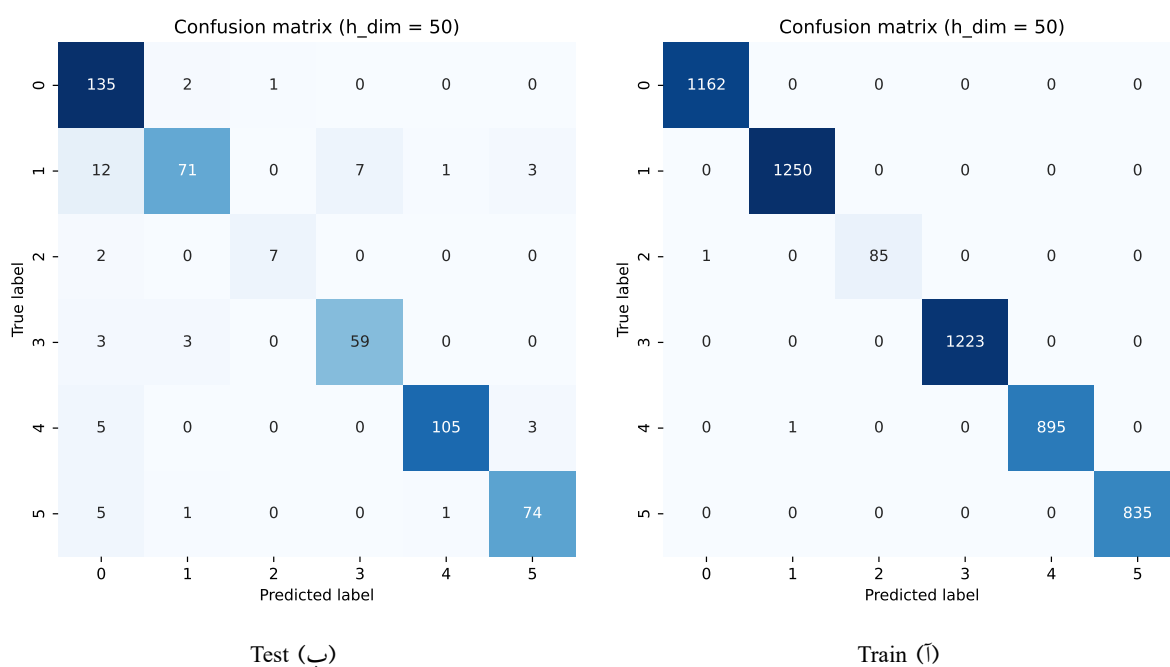


شکل ۱۰: نمودار دقت و اتلاف - مدل اول - $h=100$.

به حالت قبلی وابسته است (که هر دو با h_{t-1} و c_{t-1} نمایش داده می‌شوند) زیرا لایه نهفته‌ساز به محاسبه h_t کمکی نمی‌کند. با آموزش شبکه به صورت نظارت شده برای ارتباط زیرکلاس با این خروجی، وابستگی ای به هر دو دسته‌بندی اصلی و خود سوال ایجاد خواهد شد. دو اطلاعاتی که از بخش بازگشتی به بیرون می‌آیند، توسط دو لایه تماماًمتصل با فعال‌ساز Softmax تمیز داده و طبقه‌بندی می‌شوند؛ بنابراین، دنباله ورودی را به دو بردار با اندازه ثابت که نماینده کلاس اصلی و زیرکلاس مرتبط با آن است تبدیل می‌کنند. نمایش شماتیک مدل دوم که توصیف شد را می‌توان در **شکل ۱۵** مشاهده کرد. الگوریتم BPTT همچنان برای آموزش مدل استفاده می‌شود، اما با تفاوتی نسبت به مورد قبلی؛ که به تابع اتلاف دوگانه (به دلیل وجود دو نوع دسته‌بندی) بازمی‌گردد. برای پیاده‌سازی مدل دوم دستورات آورده‌شده در **برنامه ۲** را نوشته‌ایم. ابتدا مدل را به چند صورت پیاده‌سازی می‌کنیم. اولین و مهم‌ترین شکل پیاده‌سازی مدل که نتایج خوبی به دست داده را در این جا معرفی می‌کنیم. ابتدا لایه ورودی را

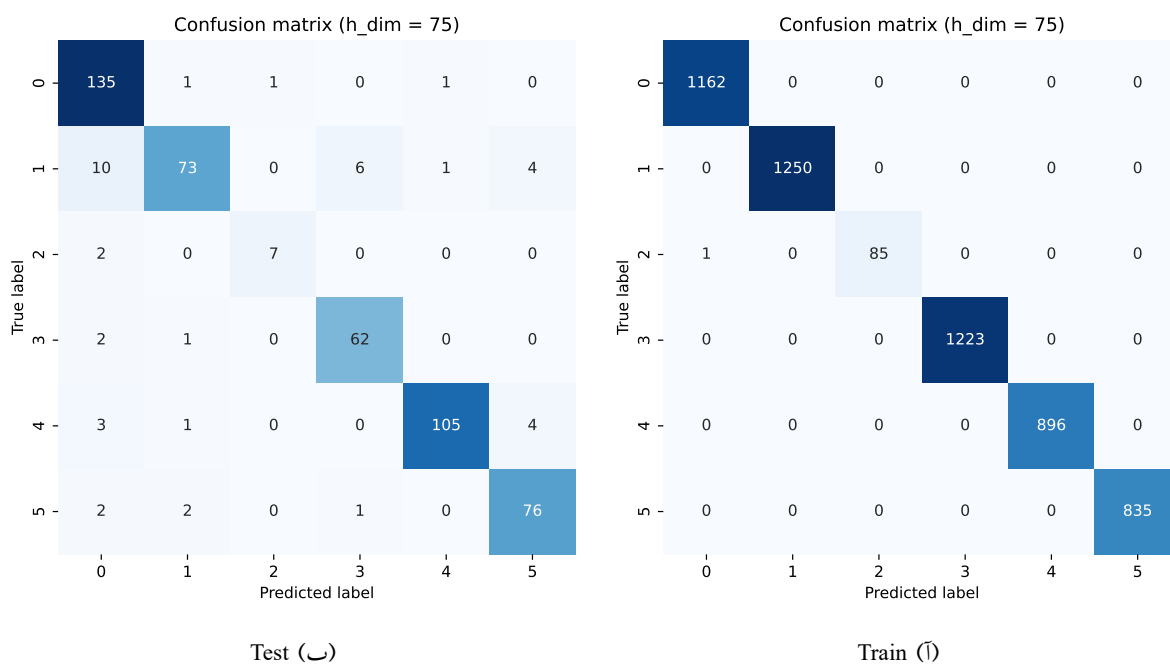


شکل ۱۱: ماتریس درهم‌ریختگی - مدل اول - h=25.

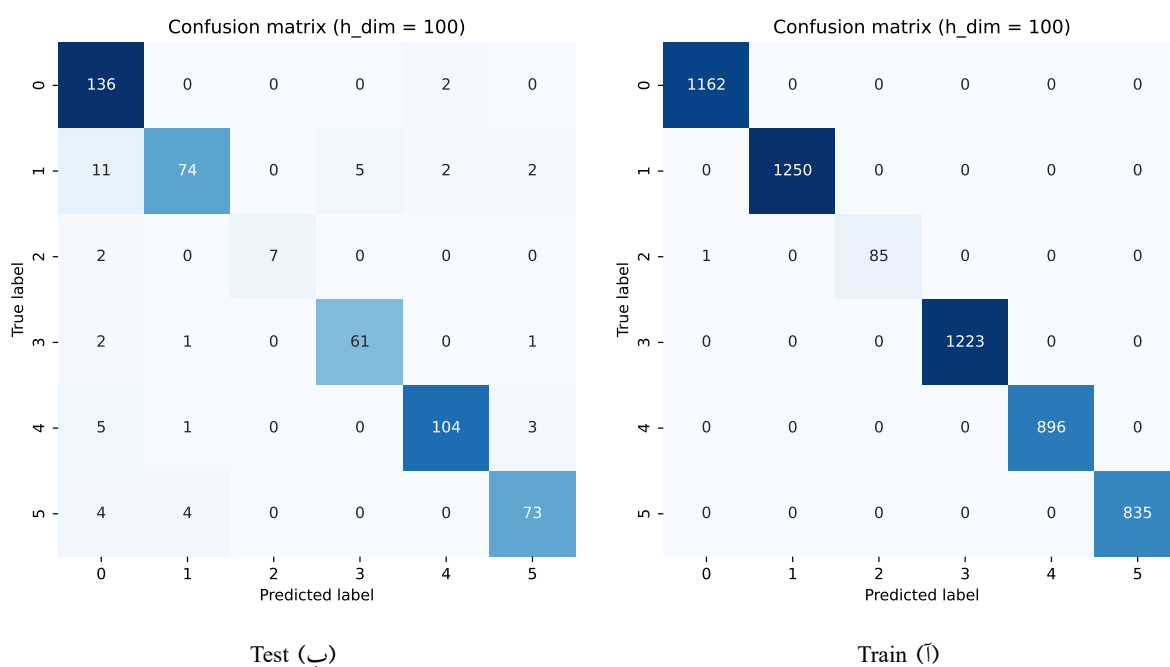


شکل ۱۲: ماتریس درهم‌ریختگی - مدل اول - h=50.

ایجاد می‌کنیم که اندازه ورودی آن با توجه به input_length تعیین می‌شود. سپس لایه Embedding را ایجاد می‌کنیم که برای تبدیل ورودی به بردارهای تعبیه استفاده می‌شود. ابعاد واژگان واحد تعبیه را با vocab_size و ابعاد بردارهای واحد تعبیه را با embedding_dim نشان می‌دهیم. ماتریس وزن‌های تعبیه embedding_matrix برای تعبیه واژگان استفاده می‌شود. این لایه در

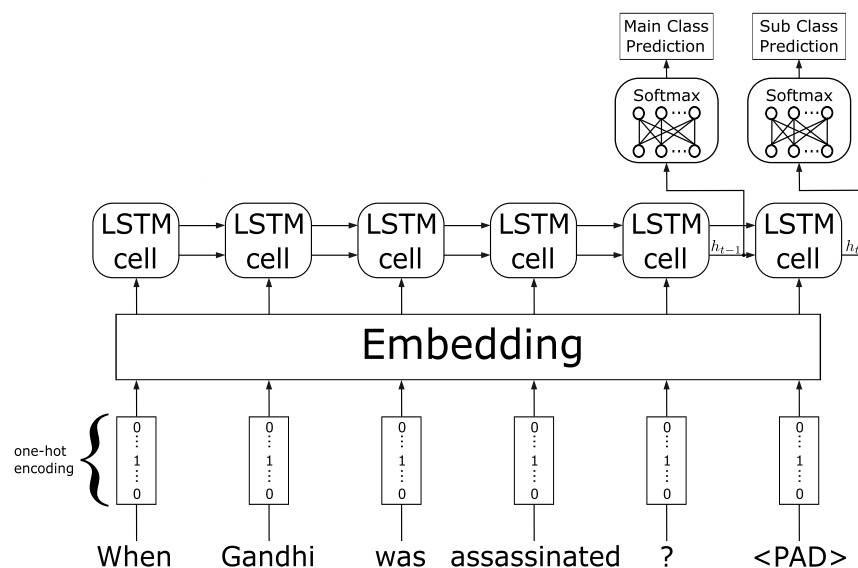


شکل ۱۳: ماتریس درهم‌ریختگی - مدل اول - h=75.



شکل ۱۴: ماتریس درهم‌ریختگی - مدل اول - h=100.

حالت غیرقابل آموزش قرار داده شده است. در ادامه یک لایه LSTM ایجاد می‌کنیم و آن را روی ورودی تعبیه‌شده اعمال می‌کنیم. تعداد واحدهای LSTM در این لایه با h_dim تعیین می‌شود. با تنظیم پارامتر return_sequences=True خروجی LSTM به صورت توالی برگردانده می‌شود و با تنظیم return_state=True وضعیت پنهان آخرین زمان برگردانده می‌شود. درواقع، این



شکل ۱۵: نمایی از مدل دوم.

لایه LSTM با دریافت ورودی دنباله‌ای (بردارهای تعبیه)، خروجی‌ها و وضعیت‌های آخرین زمان در دنباله را تولید می‌کند. نتایج این عملیات در سه متغیر `lstm_outputs`، `state_h` و `state_c` ذخیره می‌شود. متغیر `lstm_outputs` حاوی خروجی‌های LSTM برای تمام زمان‌ها در دنباله است. اگر طول دنباله T باشد، `lstm_outputs` یک تانسور با ابعاد $(batch_size, T)$ است. هر عنصر از این تانسور نشان‌دهنده خروجی LSTM در زمان مربوطه است. متغیر `state_h` نماینده حالت پنهان آخرین زمان در دنباله است. اگر h_dim اندازه بردار حالت پنهان باشد، `state_h` یک تانسور با ابعاد $(batch_size, h_dim)$ است. متغیر `state_c` حاوی وضعیت سلول آخرین زمان در دنباله و نیز یک تانسور با ابعاد $(batch_size, h_dim)$ است. در قسمتی از دستورات مقداردهی چندگانه به کمک عملگر تجزیه‌ای (=) استفاده شده است؛ بنابراین مقادیر محاسبه شده توسط لایه LSTM به ترتیب در سه متغیر `lstm_outputs`، `state_h` و `state_c` قرار می‌گیرند. در ادامه لایه‌های تماماًمتصل را تعریف می‌کنیم. تعداد واحدها را برابر با تعداد کلاس‌های موجود در متغیر `encoder_main` در نظر می‌گیریم و از فعالساز `softmax` برای تولید احتمالات دسته‌بندی استفاده می‌کنیم. خروجی این لایه برچسب‌های اصلی است. هم‌چنین یک لایه تماماًمتصل دیگر بر اساس تعداد کلاس‌های موجود در متغیر `encoder_sub` در خروجی زمان آخر LSTM تعریف می‌کنیم. بعد از تعریف لایه‌ها، مدل کلی را تعریف و مشخصات آموزش، از نوع و پارامترهای بهینه‌ساز و تابع اتلاف تا سایر موارد را به دقت تعریف می‌کنیم. در ادامه تابع `save_plots_to_pdf` را که در قسمت قبل ایجاد کرده بودیم برای نمایش نمودارهای مربوط به دقت و خطا برای دو برچسب اصلی و فرعی بازنویسی می‌کنیم. دستورات مربوط به رسم ماتریس درهم‌ریختگی هم در تابع `save_confusion_matrix_to_pdf` به صورتی مناسب بازنویسی شده‌اند که دو ماتریس درهم‌ریختگی را برای کلاس‌های یکتای موجود در برچسب‌های اصلی و فرعی نشان دهند. باید دقت داشت که برخی شماره‌های کلاس در برچسب فرعی دیتاست تست ناموجود هستند؛ مانند کلاس ۱.

```
1 Available classes in the test dataset:
2 [ 0  2  3  4  5  6  7  8  9 10 11 12 13 14 17 18 19 20 21 22 23 24 25 26
3  27 28 29 30 31 33 34 35 37 38 40 41 44 45 46]
```

در ادامه و پس از نوشتن تابعی برای محاسبه و نمایش معیارهای ارزیابی، مشابه قسمت قبل بخشی را برای آماده‌سازی داده‌ها

می‌نویسیم. برای این کار ابتدا داده‌های آموزش و آزمون به‌صورت جداگانه بارگذاری می‌شوند. سپس تابعی برای تصحیح و نرمال‌سازی متن ایجاد می‌شود. ابتدا تمام کاراکترهای غیرضروری و اعداد را حذف و با فاصله جایگزین می‌کنیم و سپس تمام حروف را به حروف کوچک تبدیل می‌کنیم. در ادامه، از کلاس `Tokenizer` برای تبدیل متن به توکن‌ها استفاده می‌کنیم. ابتدا یک نمونه از کلاس `Tokenizer` ایجاد می‌کنیم و با استفاده از تابع `fit_on_texts`، توکن‌ها از متن‌های آموزش استخراج می‌شوند و در داخل تابع ایجاد شده ذخیره می‌شوند. سپس توکن‌ها را به توالی‌ها تبدیل می‌کنیم. از تابع `texts_to_sequences` برای تبدیل متن‌های آموزش و آزمون به توالی‌های عددی استفاده می‌شود. طول بزرگترین توالی در توالی‌های آموزش را مبنا گرفته و با استفاده از تابع `pad_sequences`، توالی‌ها را به همان طول پدینگ می‌کنیم. هم‌چنین، برچسب‌های اصلی و فرعی را به‌صورت One-Hot رمزگذاری می‌کنیم. از کلاس `LabelEncoder` برای رمزگذاری برچسب‌ها استفاده می‌کنیم. در انتهای بخش مربوط به آماده‌سازی داده، برای بارگیری و استفاده از توکن‌های `GloVe` به منظور ساختن ماتریس تعبیه برای استفاده در مدل دستوراتی مشابه قبل نوشته‌ایم. در اینجا، ماتریس تعبیه نشان می‌دهد که هر توکن یا کلمه در مدل با یک بردار عددی نشان داده می‌شود. در نهایت یک حلقه می‌نویسیم تا به‌ازای پارامترهای تعیین‌شده مدل را آموزش دهد، ارزیابی کند و تمامی موارد مورد نیاز را در خروجی نمایش دهد.

Program 2: Model Two, Training, and Evaluation

```
1 import pandas as pd
2 import numpy as np
3 import re
4 import nltk
5 import tensorflow as tf
6 from nltk.tokenize import word_tokenize
7 from tensorflow.keras.preprocessing.text import Tokenizer
8 from tensorflow.keras.preprocessing.sequence import pad_sequences
9 from sklearn.preprocessing import LabelEncoder
10 from tensorflow.keras.layers import Input, Embedding, LSTM, Dense
11 from tensorflow.keras.models import Model
12 from tensorflow.keras.optimizers import Adam
13 from sklearn.metrics import f1_score, precision_score, recall_score, confusion_matrix
14 import seaborn as sns
15 import matplotlib.pyplot as plt
16 from matplotlib.backends.backend_pdf import PdfPages
17 import time
18 import warnings
19 from tensorflow.keras.utils import to_categorical
20
21 # Ignore all warnings
22 warnings.filterwarnings("ignore")
23
24 # Set random seeds for reproducibility
25 np.random.seed(42)
26 tf.random.set_seed(42)
27
28 # Set hyperparameters
29 h_dimensions = [25, 100]
```



```

30 epochs = 50
31 batch_size = 32
32
33 # Function to create the LSTM model
34 def create_lstm_model(input_length, h_dim):
35     input_layer = Input(shape=(input_length,))
36     embedding_layer = Embedding(vocab_size, embedding_dim, weights=[embedding_matrix], trainable=
False)(input_layer)
37     lstm_layer = LSTM(h_dim, activation='tanh', return_sequences=True, return_state=True)
38     lstm_outputs, state_h, state_c = lstm_layer(embedding_layer)
39     main_output = Dense(len(encoder_main.classes_), activation='softmax', name='main_output')(
state_h)
40     sub_output = Dense(len(encoder_sub.classes_), activation='softmax', name='sub_output')(
lstm_outputs[:, -1, :])
41     model = Model(inputs=input_layer, outputs=[main_output, sub_output])
42     model.compile(optimizer=Adam(learning_rate=0.001), loss={'main_output': '
sparse_categorical_crossentropy', 'sub_output': 'sparse_categorical_crossentropy'}, metrics=[
'accuracy'])
43     return model
44
45 # Function to save plots to PDF
46 def save_plots_to_pdf(filename, history, h_dim):
47     with PdfPages(filename) as pdf:
48         fig, axs = plt.subplots(2, 2, figsize=(10, 8))
49
50         axs[0, 0].plot(history.history['main_output_accuracy'], label='Training')
51         axs[0, 0].plot(history.history['val_main_output_accuracy'], label='Validation')
52         axs[0, 0].set_xlabel('Epoch')
53         axs[0, 0].set_ylabel('Accuracy')
54         axs[0, 0].legend()
55         axs[0, 0].set_title(f'Main Accuracy (h_dim = {h_dim})')
56
57         axs[0, 1].plot(history.history['main_output_loss'], label='Training')
58         axs[0, 1].plot(history.history['val_main_output_loss'], label='Validation')
59         axs[0, 1].set_xlabel('Epoch')
60         axs[0, 1].set_ylabel('Loss')
61         axs[0, 1].legend()
62         axs[0, 1].set_title(f'Main Loss (h_dim = {h_dim})')
63
64         axs[1, 0].plot(history.history['sub_output_accuracy'], label='Training')
65         axs[1, 0].plot(history.history['val_sub_output_accuracy'], label='Validation')
66         axs[1, 0].set_xlabel('Epoch')
67         axs[1, 0].set_ylabel('Accuracy')
68         axs[1, 0].legend()
69         axs[1, 0].set_title(f'Sub Accuracy (h_dim = {h_dim})')

```

```

70
71     axs[1, 1].plot(history.history['sub_output_loss'], label='Training')
72     axs[1, 1].plot(history.history['val_sub_output_loss'], label='Validation')
73     axs[1, 1].set_xlabel('Epoch')
74     axs[1, 1].set_ylabel('Loss')
75     axs[1, 1].legend()
76     axs[1, 1].set_title(f'Sub Loss (h_dim = {h_dim})')
77
78     plt.subplots_adjust(hspace=0.5) # Adjust the vertical spacing between subplots
79
80     pdf.savefig(bbox_inches='tight')
81
82     # Show the plots
83     plt.show()
84
85     plt.close()
86
87 # Function to save confusion matrix to PDF
88 def save_confusion_matrix_to_pdf(y_true, y_pred, filename, label_type):
89     unique_classes = np.unique(y_true)
90     num_classes = len(unique_classes)
91
92     cm = confusion_matrix(y_true, y_pred, labels=unique_classes)
93     cm_normalized = cm.astype('float') / cm.sum(axis=1)[:, np.newaxis] # Normalize the confusion
94     matrix
95
96     # Calculate the appropriate figure size based on the number of classes
97     figsize = max(10, num_classes / 2)
98
99     plt.figure(figsize=(figsize, figsize))
100     sns.heatmap(cm_normalized, annot=True, fmt='.2f', cmap='Blues', cbar=False)
101
102     # Adjust the font size of the numbers inside the heatmap
103     annot_fontsize = max(6, 120 // num_classes)
104     plt.tick_params(axis='both', labelsize=annot_fontsize)
105
106     # Map the class indices to their corresponding labels
107     class_labels = encoder_sub.inverse_transform(unique_classes)
108
109     plt.xticks(np.arange(num_classes) + 0.5, class_labels, rotation='vertical')
110     plt.yticks(np.arange(num_classes) + 0.5, class_labels, rotation='horizontal')
111
112     plt.xlabel('Predicted label')
113     plt.ylabel('True label')
114     plt.title(f'Normalized Confusion Matrix ({label_type})')

```

```

114 plt.savefig(filename, bbox_inches='tight')
115
116 # Show the plot
117 plt.show()
118 plt.close()
119
120
121 # Print the available classes in the test dataset
122 print("Available classes in the test dataset:")
123 print(np.unique(test_encoded_sub_labels))
124
125 # Function to print evaluation metrics
126 def print_evaluation_metrics(y_true, y_pred, label_type):
127     accuracy = np.mean(y_true == y_pred)
128     precision = precision_score(y_true, y_pred, average='macro')
129     recall = recall_score(y_true, y_pred, average='macro')
130     f1 = f1_score(y_true, y_pred, average='macro')
131
132     print(f'Evaluation metrics ({label_type}):')
133     print(f'Accuracy: {accuracy:.4f}')
134     print(f'Precision: {precision:.4f}')
135     print(f'Recall: {recall:.4f}')
136     print(f'F1-score: {f1:.4f}')
137
138 # Load dataset
139 train_data = pd.read_csv('train.csv')
140 test_data = pd.read_csv('test.csv')
141
142 # Define function for text normalization
143 def normalize_text(text):
144     text = re.sub(r'[^a-zA-Z0-9\?]+', ' ', text)
145     text = text.lower()
146     return text
147
148 # Normalize text
149 train_data['text'] = train_data['text'].apply(normalize_text)
150 test_data['text'] = test_data['text'].apply(normalize_text)
151
152 # Tokenize the text
153 tokenizer = Tokenizer()
154 tokenizer.fit_on_texts(train_data['text'])
155
156 # Convert tokens to sequences
157 train_sequences = tokenizer.texts_to_sequences(train_data['text'])
158 test_sequences = tokenizer.texts_to_sequences(test_data['text'])

```

```

159
160 # Pad sequences
161 maxlen = max(len(seq) for seq in train_sequences)
162 train_padded_sequences = pad_sequences(train_sequences, maxlen=maxlen)
163 test_padded_sequences = pad_sequences(test_sequences, maxlen=maxlen)
164
165 # Encode main labels
166 encoder_main = LabelEncoder()
167 encoder_main.fit(train_data['label-coarse'])
168 train_encoded_main_labels = encoder_main.transform(train_data['label-coarse'])
169 test_encoded_main_labels = encoder_main.transform(test_data['label-coarse'])
170
171 # Encode sub labels
172 encoder_sub = LabelEncoder()
173 encoder_sub.fit(train_data['label-fine'])
174 train_encoded_sub_labels = encoder_sub.transform(train_data['label-fine'])
175 test_encoded_sub_labels = encoder_sub.transform(test_data['label-fine'])
176
177 # Convert main labels to one-hot vectors
178 y_main_train = to_categorical(train_encoded_main_labels)
179 y_main_test = to_categorical(test_encoded_main_labels)
180
181 # Convert sub labels to one-hot vectors
182 y_sub_train = to_categorical(train_encoded_sub_labels)
183 y_sub_test = to_categorical(test_encoded_sub_labels)
184
185 # Load GloVe embeddings into a dictionary
186 embeddings_index = {}
187 with open('glove.6B.300d.txt') as f:
188     for line in f:
189         values = line.split()
190         word = values[0]
191         coefs = np.asarray(values[1:], dtype='float32')
192         embeddings_index[word] = coefs
193
194 # Create an embedding matrix
195 embedding_dim = 300
196 vocab_size = len(tokenizer.word_index) + 1
197 embedding_matrix = np.zeros((vocab_size, embedding_dim))
198 for word, i in tokenizer.word_index.items():
199     embedding_vector = embeddings_index.get(word)
200     if embedding_vector is not None:
201         embedding_matrix[i] = embedding_vector
202
203 # Iterate over h_dimensions and train the models

```

```

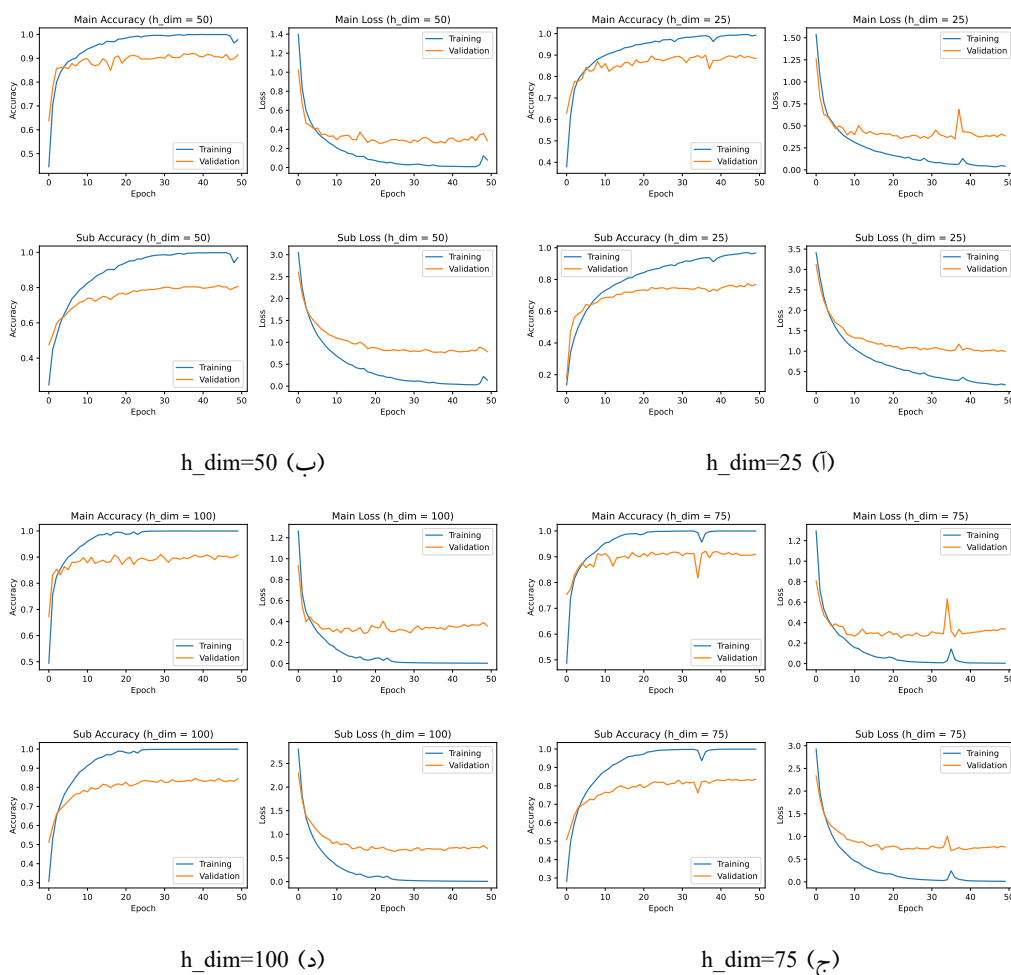
204 for h_dim in h_dimensions:
205     print(f"Training model with h_dim = {h_dim}")
206     model = create_lstm_model(maxlen, h_dim)
207
208     # Fit the model
209     history = model.fit(
210         train_padded_sequences,
211         {'main_output': train_encoded_main_labels, 'sub_output': train_encoded_sub_labels},
212         validation_data=(test_padded_sequences, {'main_output': test_encoded_main_labels, '
213         sub_output': test_encoded_sub_labels}),
214         epochs=epochs,
215         batch_size=batch_size,
216         verbose=1
217     )
218
219     # Save plots and confusion matrices
220     timestamp = str(int(time.time()))
221     save_plots_to_pdf(f'plots_{h_dim}_{timestamp}.pdf', history, h_dim)
222
223     # Only plot confusion matrices for the test data
224     main_output_predictions, sub_output_predictions = model.predict(test_padded_sequences)
225     main_output_predictions = np.argmax(main_output_predictions, axis=1)
226     sub_output_predictions = np.argmax(sub_output_predictions, axis=1)
227
228     save_confusion_matrix_to_pdf(
229         test_encoded_main_labels,
230         main_output_predictions,
231         f'confusion_matrix_test_main_{h_dim}_{timestamp}.pdf',
232         'main'
233     )
234
235     save_confusion_matrix_to_pdf(
236         test_encoded_sub_labels,
237         sub_output_predictions,
238         f'confusion_matrix_test_sub_{h_dim}_{timestamp}.pdf',
239         'sub'
240     )
241
242     print_evaluation_metrics(test_encoded_main_labels, main_output_predictions, 'Main')
243     print_evaluation_metrics(test_encoded_sub_labels, sub_output_predictions, 'Sub')

```

با اجرای این دستورات در نزدیک‌ترین حالت به مقاله نتایج به صورتی است که در جدول ۲ و شکل ۱۶ و شکل ۱۷ نشان داده شده است. همان‌طور که مشاهده می‌شود نتایج در حالت $h_dim=100$ به مراتب بهتر است. دلیل بدشدن نتایج برای کلاس و برپسی مانند ۶ یا ۲۹ فرعی آن است که فقط یک داده از این برچسب در مجموعه داده‌ی تست وجود دارد و رخ دادن هم‌چنین اتفاقی دور از انتظار نیست. عملکرد کلی اما قابل قبول است.

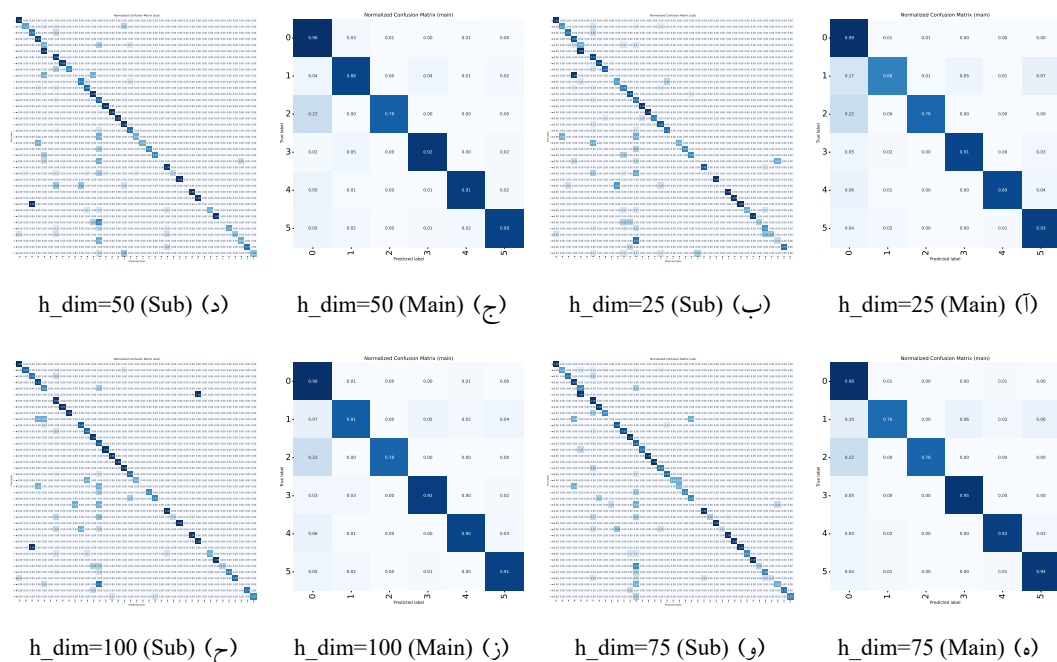
جدول ۲: نتایج مربوط به مدل دوم - پیاده‌سازی اول

Test Set (%)								Training Set (%)		h Dimension
F1-score		Recall		Precision		Accuracy		Accuracy		
Sub	Main	Sub	Main	Sub	Main	Sub	Main	Sub	Main	
۶۱.۶۲	۸۶.۵۴	۶۲.۴۸	۸۶.۱۹	۶۶.۴۵	۸۷.۹۱	۷۶.۸۰	۸۸.۴۰	۹۶.۷۴	۹۹.۲۸	25
۶۰.۷۲	۸۹.۹۱	۶۱.۱۷	۸۹.۰۱	۶۶.۱۶	۹۱.۰۰	۸۰.۶۰	۹۱.۴۰	۹۶.۹۴	۹۷.۷۶	50
۶۹.۲۶	۹۰.۳۱	۶۷.۶۵	۸۸.۷۳	۷۸.۰۱	۹۲.۸۷	۸۳.۶۰	۹۱.۰۰	۹۹.۹۳	۹۹.۹۶	75
۶۳.۱۲	۹۰.۱۶	۶۵.۶۵	۸۸.۴۰	۶۵.۹۲	۹۲.۶۷	۸۴.۴۰	۹۰.۸۰	۹۹.۹۳	۹۹.۹۶	100



شکل ۱۶: نتیجه تابع دقت و اتلاف در حالت اعتبارسنجی با داده‌های تست (منطبق بر مقاله).

نتایج را در حالتی که در لایه‌های تماماً متصل از رگولاریزیشن نرم L2 استفاده کرده‌ایم نیز تست کردیم. نتایج به صورتی است که در جدول ۳ و شکل ۱۸ و شکل ۱۹ نشان داده شده است. هم‌چنین نتایج را برای یک حالتی که احتمالاً با آنچه مقاله در خصوص داده‌های اعتبارسنجی در آن به‌صورت گنگ اشاره کرده کمی فاصله دارد نیز امتحان کردیم (انتخاب بخشی از داده آموزش به‌عنوان اعتبارسنجی). نتایج به صورتی است که در جدول ۴



شکل ۱۷: ماتریس درهم ریختگی در حالت اعتبارسنجی با داده های تست (منطبق بر مقاله).

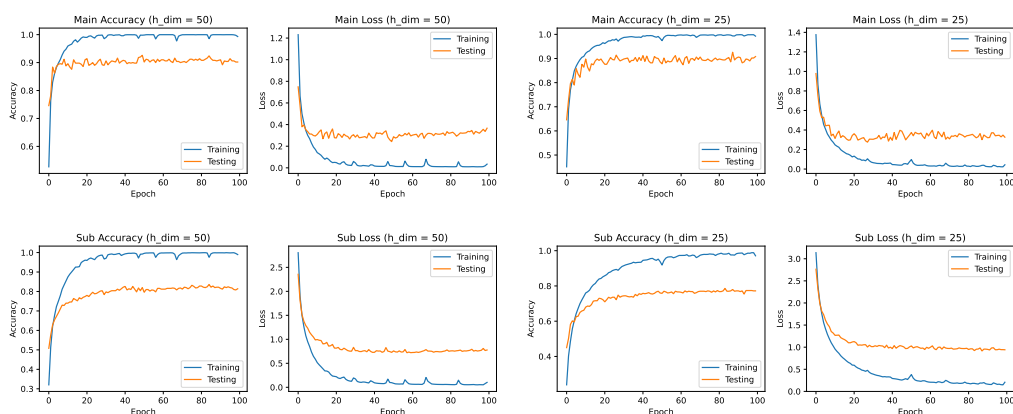
جدول ۳: نتایج مربوط به مدل دوم - پیاده سازی دوم

Test Set (%)								Training Set (%)		h Dimension
F1-score		Recall		Precision		Accuracy		Accuracy		
Sub	Main	Sub	Main	Sub	Main	Sub	Main	Sub	Main	
۵۰.۶۶	۸۹.۱۲	۵۲.۴۳	۸۸.۵۶	۵۳.۹۸	۸۹.۸۸	۷۷.۲۰	۹۰.۸۰	۹۷.۱۰	۹۹.۲۷	25
۶۸.۱۲	۸۸.۷۸	۶۷.۶۲	۸۶.۴۳	۷۷.۶۸	۹۳.۲۰	۸۱.۴۰	۹۰.۲۰	۹۸.۹۷	۹۹.۳۲	50
۶۴.۸۷	۸۹.۴	۶۵.۴۶	۸۷.۴۵	۶۸.۲۴	۹۲.۴۳	۸۲.۸۰	۸۹.۸۰	۹۹.۹۳	۹۹.۹۸	75
۶۷.۰۴	۹۰.۰۸	۶۶.۳۸	۸۸.۷۶	۷۱.۹۵	۹۳.۱۰	۸۲.۴۰	۹۱.۲۰	۹۹.۴۹	۹۹.۷۴	100

و شکل ۲۰ و شکل ۲۱ نشان داده شده است. در نهایت، حالت قبل را با ترم تنظیم تکرار می کنیم. نتایج به صورتی است که در

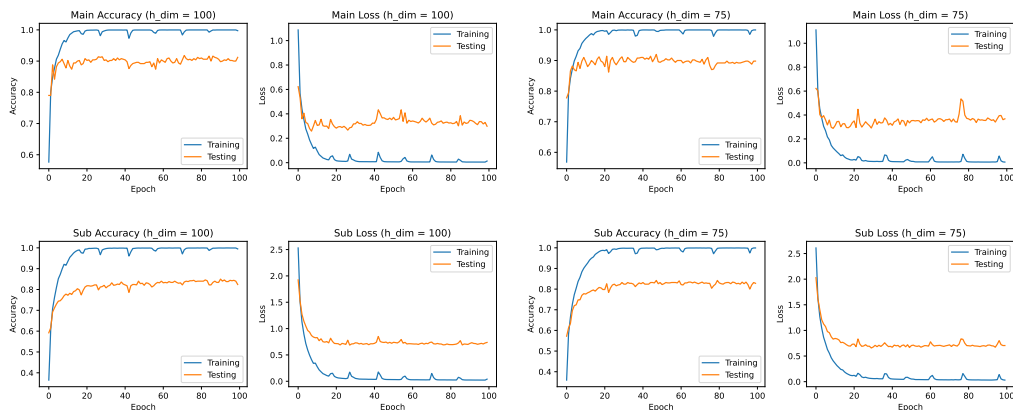
جدول ۴: نتایج مربوط به مدل دوم - پیاده سازی سوم

Test Set (%)								Training Set (%)		h Dimension
F1-score		Recall		Precision		Accuracy		Accuracy		
Sub	Main	Sub	Main	Sub	Main	Sub	Main	Sub	Main	
۴۶.۲۹	۸۷.۸۳	۴۷.۴۹	۸۴.۱۷	۵۱.۲۱	۸۶.۵۹	۷۵.۲۰	۸۷.۶۰	۹۵.۰۵	۹۹.۳۸	25
۶۵.۵۲	۸۸.۲۶	۶۴.۹۳	۸۶.۴۶	۷۰.۳۰	۹۱.۱۳	۸۱.۲۰	۸۸.۸۰	۹۹.۸۲	۹۹.۹۳	50
۶۲.۴۴	۸۸.۷۹	۶۰.۸۲	۸۶.۶۸	۶۸.۵۱	۹۲.۲۱	۷۹.۲۰	۸۹.۰۰	۹۹.۸۹	۹۹.۹۵	75
۶۳.۷۳	۸۸.۸۹	۶۳.۴۸	۸۷.۷۶	۶۹.۱۴	۹۰.۸۹	۸۲.۴۰	۹۰.۲۰	۹۹.۹۳	۹۹.۹۵	100



h_dim=50 (ب)

h_dim=25 (ا)



h_dim=100 (د)

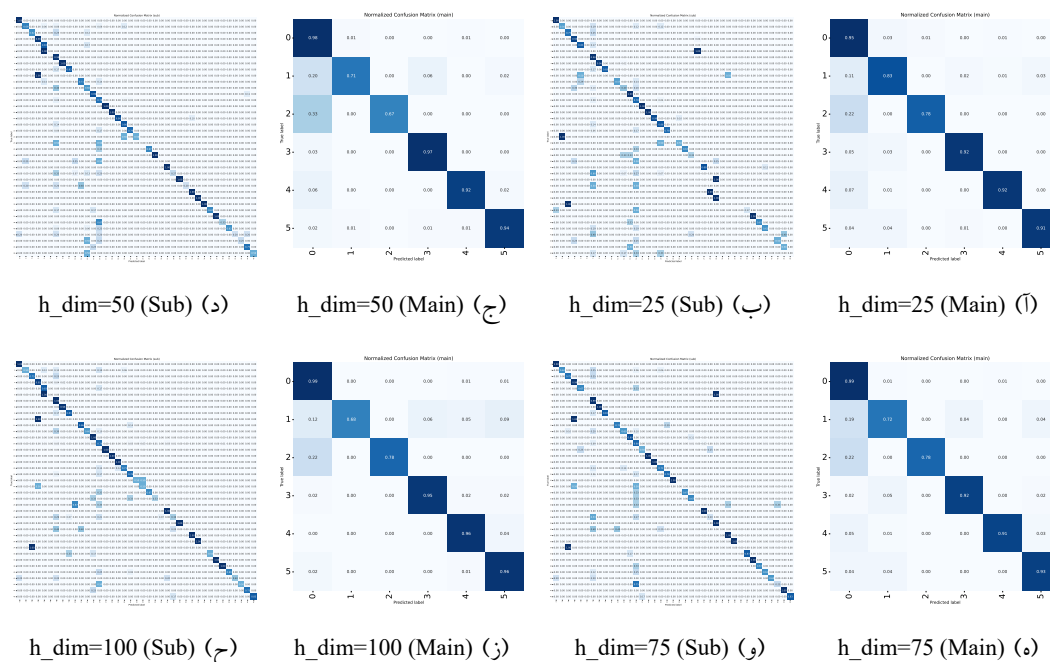
h_dim=75 (ج)

شکل ۱۸: نتیجه تابع دقت و اتلاف در حالت اعتبارسنجی با داده‌های تست (با ترم تنظیم).

جدول ۵ و شکل ۲۲ و شکل ۲۳ نشان داده شده است.

جدول ۵: نتایج مربوط به مدل دوم - پیاده‌سازی چهارم

Test Set (%)								Training Set (%)		h Dimension
F1-score		Recall		Precision		Accuracy		Accuracy		
Sub	Main	Sub	Main	Sub	Main	Sub	Main	Sub	Main	
۴۸.۵۵	۸۸.۴۴	۵۰.۱۵	۸۶.۴۰	۵۴.۴۵	۹۱.۶۹	۷۶.۲۰	۸۸.۸۰	۹۷.۶۶	۹۹.۸۲	25
۶۵.۸۰	۹۰.۰۴	۶۵.۹۱	۸۸.۰۴	۷۰.۸۴	۹۲.۹۸	۸۱.۰۰	۹۰.۴۰	۹۹.۸۹	۹۹.۹۸	50
۵۹.۹۸	۸۸.۳۷	۶۰.۰۷	۸۶.۸۶	۶۴.۲۲	۹۰.۶۴	۷۸.۰۰	۸۹.۴۰	۹۸.۳۰	۹۹.۱۵	75
۶۴.۹۸	۸۷.۴۲	۶۴.۰۸	۸۵.۰۲	۷۱.۴۳	۹۱.۶۸	۸۰.۸۰	۸۹.۰۰	۹۹.۳۸	۹۹.۵۹	100

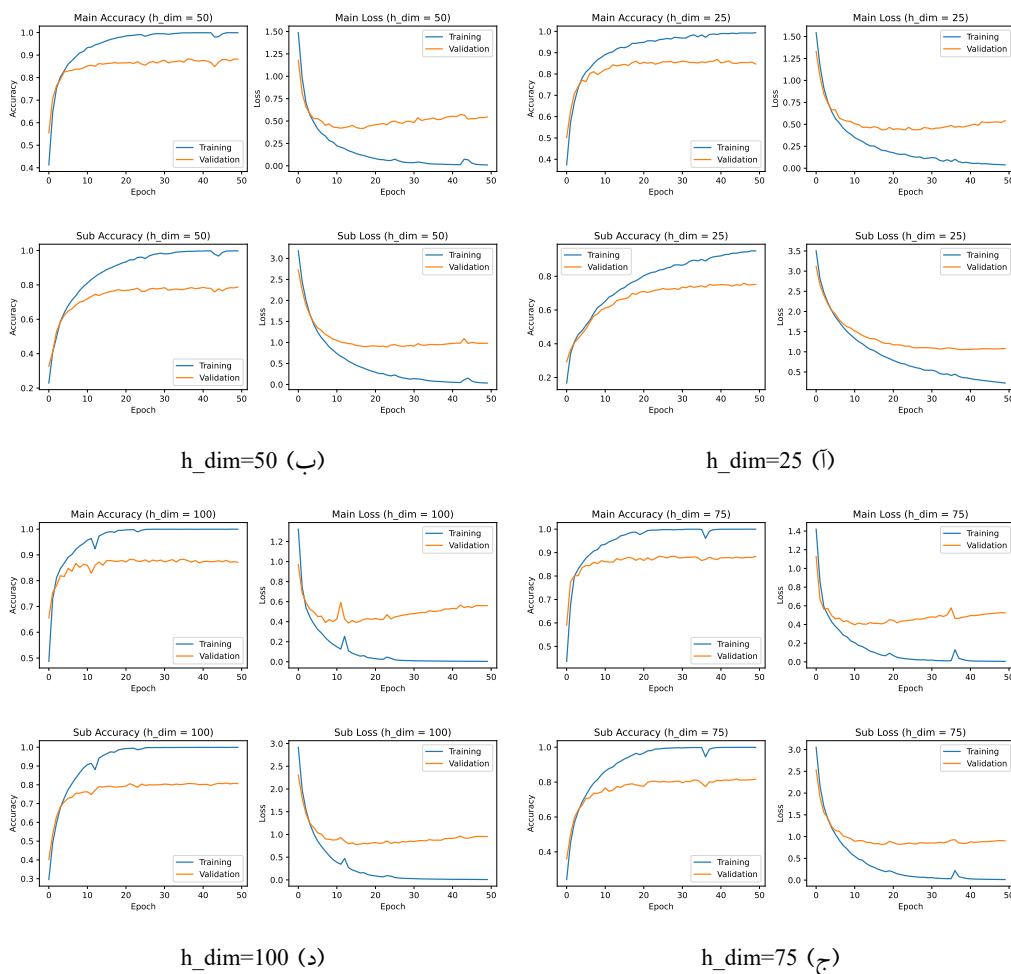


شکل ۱۹: ماتریس درهم‌ریختگی در حالت اعتبارسنجی با داده‌های تست (با ترم تنظیم).

۴.۱ پاسخ قسمت ۴

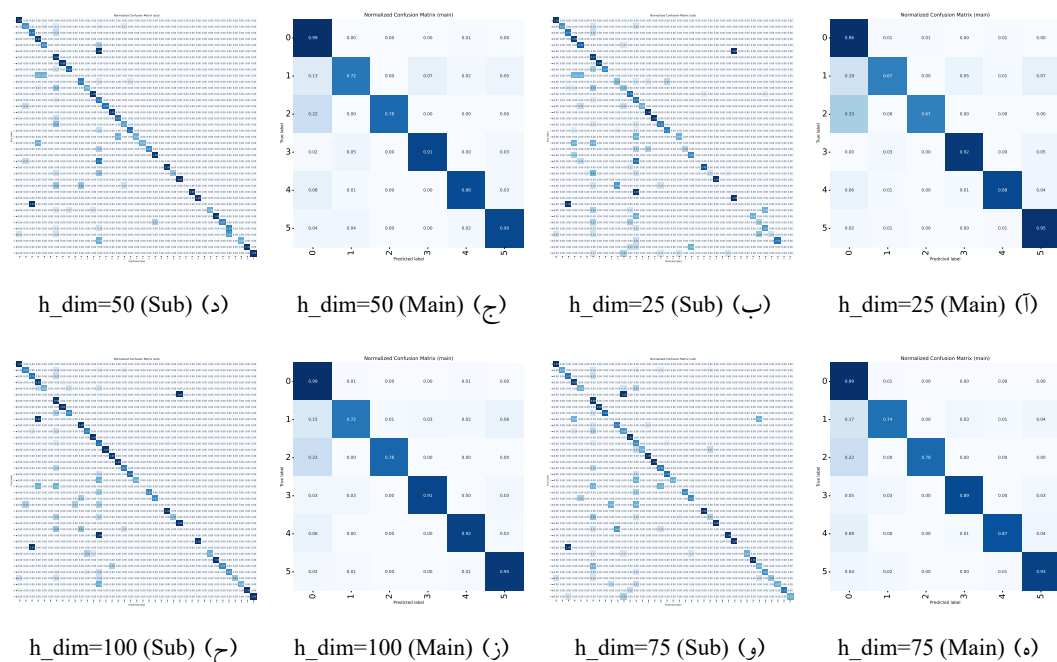
در مقاله این‌طور ذکر شده که در نهایت و به منظور آزمایش (حداقل به طور خلاصه) اهمیت آنچه در طبقه‌بندی اهداف به منظور خلق یک پاسخ‌دهنده به دست آمده، تصمیم گرفته شد یک نمونه اولیه از یک مدل پاسخ‌دهنده ایجاد شود که بتواند از ساختار به دست آمده استفاده کند و یک پاسخ تولید کند. این نمونه اولیه از یک شبکه LSTM دوطرفه (BLSTM) استفاده می‌کند تا بتواند بررسی ورودی را به صورت مناسبی انجام دهد، به طوری که تنها پس از دریافت کامل سوال، بتواند پاسخی تولید کند. وضعیت این شبکه BLSTM توسط پیش‌بینی بازگشتی توسط مدل قبلی تغذیه می‌شود. دو بردار که نماینده کلاس اصلی و فرعی هستند در واقع به هم متصل شده‌اند و یک بردار تکی را تشکیل می‌دهند که وضعیت اولیه شبکه BLSTM را مشخص کنند. خروجی‌های شبکه، که کلمات پاسخ را تشکیل می‌دهند، به عبارت دیگر توسط تحلیل هر دو متن قبلی و بعدی تولید می‌شوند و تحت تأثیر دسته‌بندی فراهم‌شده قرار می‌گیرند. نمایی از این مدل پاسخ‌دهنده (پاسخ‌گو) در شکل ۲۵ نشان داده شده است.

برای تشکیل ساختار آموزش و ارزیابی این بخش هم مانند قبل ابتدا بخشی را برای پردازش و آماده‌سازی داده‌ها در نظر می‌گیریم. ابتدا فایل با فرمت CSV به نام QA_data.csv را بارگیری می‌کنیم و آن را در متغیر qa_data ذخیره می‌کنیم. پارامتر encoding=latin1 برای تعیین کدگذاری متن استفاده می‌شود که در اینجا از latin1 استفاده شده است. در ادامه یک تابع با نام normalize_text تعریف می‌کنیم. این تابع وظیفه نرمال‌سازی یک متن را بر عهده دارد و ضمن کوچک کردن حروف، واژگان و نشانه‌های اضافی را حذف می‌کند. در ادامه سایر کارهای مربوط به توکنی‌سازی، نهفته‌سازی، پدگذاری و غیره مشابه آنچه که در پاسخ قسمت ۳ ذکر شد انجام می‌شود. در ادامه مدل را مطابق آنچه از مقاله فهمیده‌ایم تعریف و پیاده‌سازی می‌کنیم. مشابه قسمت قبل چند مدل مختلف ایجاد کرده‌ایم که نزدیک‌ترین مدل به مقاله را این‌جا معرفی می‌کنیم. ابتدا یک لایه ورودی با شکل مشخص که به عنوان ورودی دنباله‌های توکن‌شده با طول بیشینه max_sequence_length را دریافت می‌کند، ایجاد می‌کنیم. سپس یک لایه embedding ایجاد می‌کنیم که مسئول نمایش عددی لغات است. با ارائه پارامترهای vocab_size و embedding_dim، لایه embedding با سایز و ابعاد مشخص ایجاد می‌شود. همچنین با ارائه اوزان embedding_matrix



شکل ۲۰: نتیجه تابع دقت و اتلاف در حالت داده‌های اعتبارسنجی جداگانه.

به‌عنوان ورودی، لایه embedding از این وزن‌ها استفاده می‌کند. پارامتر trainable=False نیز به معنی این است که وزن‌های این لایه قابل آموزش نباشند. در ادامه یک لایه LSTM با ۱۰۰ واحد حافظه memory (units) ایجاد می‌کنیم که دنباله‌های ورودی را دریافت می‌کند و دنباله خروجی را ضمن داشتن پارامتر return_sequences=True تولید می‌کند. آخرین وضعیت مخفی (hidden state) و آخرین وضعیت سلول (cell state) لایه LSTM را استخراج می‌کنیم. این وضعیت‌ها معمولاً در زمان پایان دنباله معنادار هستند. سپس یک لایه BiLSTM (دوطرفه) ایجاد می‌کنیم. این لایه دو لایه LSTM را به صورت موازی در دو جهت (جلو و عقب) اجرا می‌کند و خروجی را تولید می‌کند. پارامتر merge_mode=concat نیز نشان‌دهنده این است که خروجی‌های دو لایه LSTM دوطرفه به صورت ترکیب‌شده (concatenated) باشند. این لایه از خروجی لایه embedding و بازگشتی می‌تواند استفاده کند و حالات مختلف در پیاده‌سازی تست شده‌اند. در بخش وضعیت اولیه (initial_state)، هم، آرایه‌ای شامل آخرین وضعیت مخفی و آخرین وضعیت سلول لایه LSTM آمده است (منطبق بر مقاله). این وضعیت‌ها از لایه LSTM قبلی استخراج شده و به عنوان وضعیت اولیه برای لایه BiLSTM مورد استفاده قرار می‌گیرند. در نهایت، یک لایه خروجی ایجاد می‌کنیم که دارای تابع فعال‌سازی softmax است که برای تبدیل خروجی‌های عددی به احتمالات برای هر کلاس استفاده می‌شود. تعداد کلاس‌ها با استفاده از پارامتر num_classes مشخص شده است. در انتها مدل را آماده‌به‌کار می‌کنیم. در ادامه دستوراتی را برای آموزش مدل می‌نویسیم. با فراخوانی تابع fit روی مدل، داده‌های آموزش



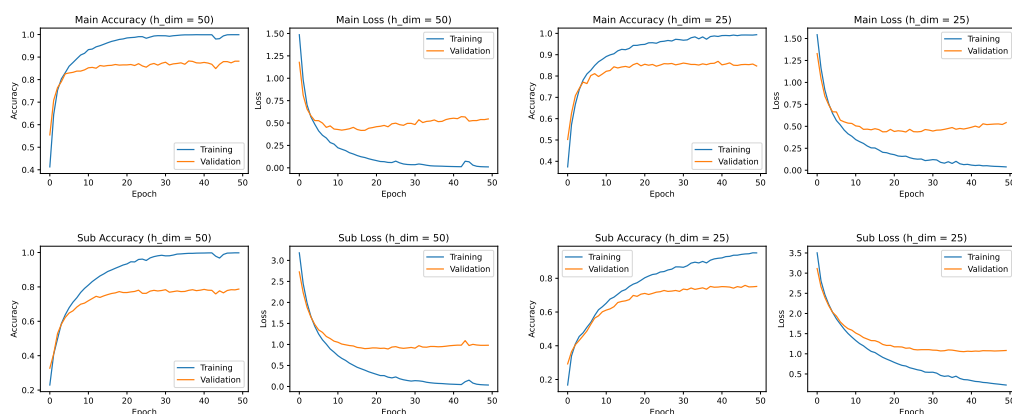
شکل ۲۱: ماتریس درهم‌ریختگی در حالت داده‌های اعتبارسنجی جداگانه.

X و برچسب‌های آموزش answers_encoded به مدل داده می‌شوند. فرآپارامترها با سعی و خطای فراوان تنظیم شده‌اند. در ادامه لیستی از سوالات مدنظر تمرین ایجاد می‌کنیم و این سوالات آزمایشی را به صورت نرمال شده (توسط تابع normalize_text) در آرایه test_questions ذخیره می‌کنیم. این عمل باعث حذف نویزهای غیرضروری و استانداردسازی متن سوالات می‌شود. در انتها سوالات آزمایشی به توکن‌ها تبدیل می‌شوند.

در ادامه ابتدا تابع texts_to_sequences از توکن‌بندی‌کننده tokenizer روی سوالات آزمایشی اعمال می‌شود و سپس با استفاده از تابع pad_sequences، دنباله‌های توکن‌ها با طول یکسان (max_sequence_length) تنظیم می‌شوند. نتیجه نهایی در آرایه test_X قرار می‌گیرد. با انجام این کار و با استفاده از مدل آموزش دیده، پیش‌بینی‌های مدل روی داده‌های آزمایشی صورت می‌گیرد. تابع predict روی test_X اجرا شده و خروجی پیش‌بینی‌ها در متغیر predictions قرار می‌گیرد. سپس با استفاده از np.argmax، برچسب‌هایی که احتمال بالاتری دارند انتخاب می‌شوند و در آرایه predicted_labels ذخیره می‌شوند. در نهایت، با استفاده از حلقه for، سوالات آزمایشی و پاسخ‌های پیش‌بینی شده به صورت همزمان چاپ می‌شوند. دستورات در برنامه ۳ آورده شده است.

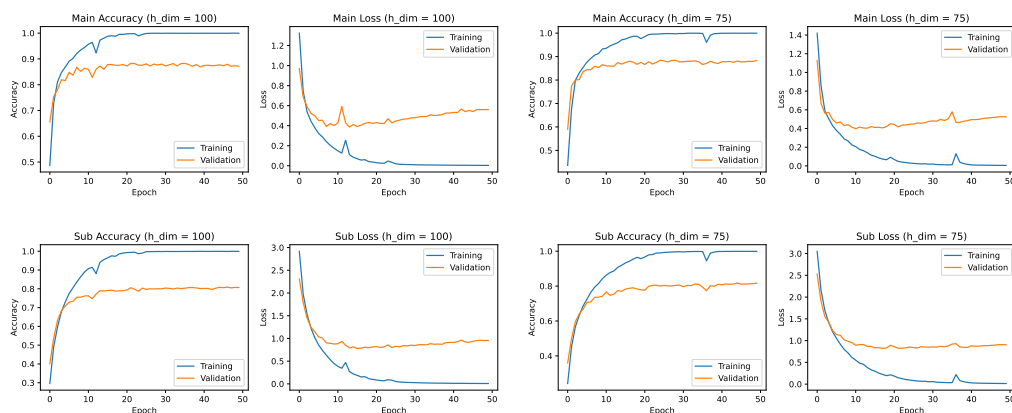
Program 3: Prototype Responder, Training, and Evaluation

```
1 import pandas as pd
2 import numpy as np
3 import re
4 import matplotlib.pyplot as plt
5 from tensorflow.keras.preprocessing.text import Tokenizer
6 from tensorflow.keras.preprocessing.sequence import pad_sequences
7 from tensorflow.keras.models import Model
8 from tensorflow.keras.layers import Input, Embedding, LSTM, Bidirectional, Dense
9
```



h_dim=50 (ب)

h_dim=25 (ا)



h_dim=100 (د)

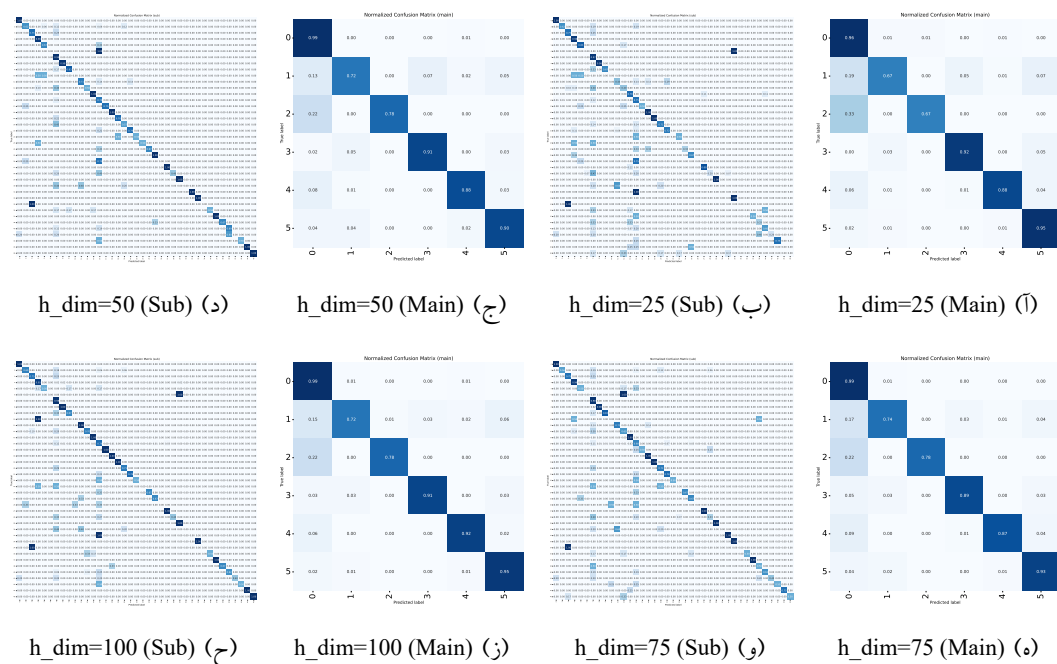
h_dim=75 (ج)

شکل ۲۲: نتیجه تابع دقت و اتلاف در حالت داده‌های اعتبارسنجی جداگانه.

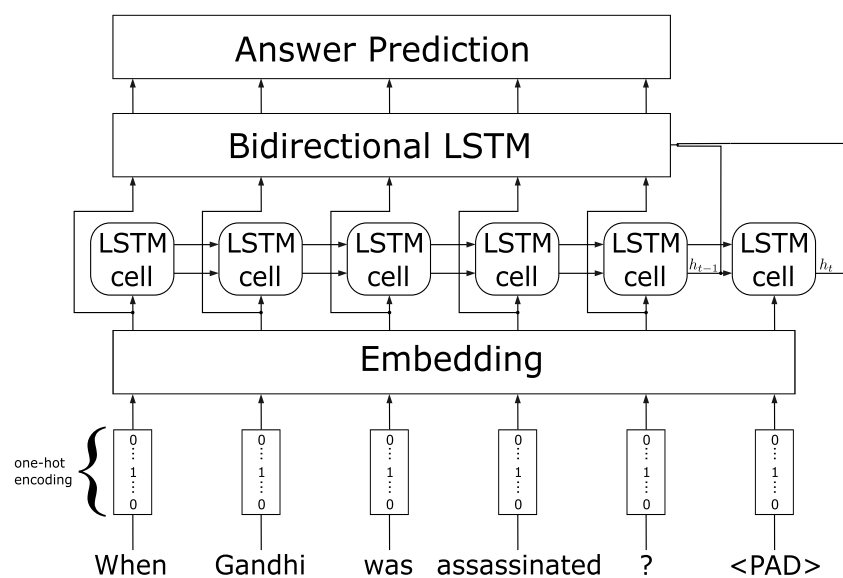
```

10 # Load the QA_data
11 qa_data = pd.read_csv('QA_data.csv', encoding='latin1')
12
13 # Define function for text normalization
14 def normalize_text(text):
15     text = re.sub(r'[^a-zA-Z0-9\?]+', ' ', text)
16     text = text.lower()
17     return text
18
19 # Normalize text
20 qa_data['text'] = qa_data['text'].apply(normalize_text)
21
22 # Preprocess the data
23 tokenizer = Tokenizer()
24 tokenizer.fit_on_texts(qa_data['text'])
25 sequences = tokenizer.texts_to_sequences(qa_data['text'])

```



شکل ۲۳: ماتریس درهم‌ریختگی در حالت داده‌های اعتبارسنجی جداگانه.



شکل ۲۴: نمایی از مدل پاسخگو.

```

26 word_index = tokenizer.word_index
27 max_sequence_length = max(len(seq) for seq in sequences)
28
29 X = pad_sequences(sequences, maxlen=max_sequence_length)
30

```

```

31 # One-hot encode the answers
32 answers = qa_data['answer']
33 answer_labels = np.unique(answers)
34 label_to_index = {label: index for index, label in enumerate(answer_labels)}
35 answers_encoded = np.array([label_to_index[answer] for answer in answers])
36
37 num_classes = len(answer_labels)
38
39 # Load GloVe embeddings into a dictionary
40 embeddings_index = {}
41 with open('glove.6B.300d.txt') as f:
42     for line in f:
43         values = line.split()
44         word = values[0]
45         coefs = np.asarray(values[1:], dtype='float32')
46         embeddings_index[word] = coefs
47
48 # Create an embedding matrix
49 embedding_dim = 300
50 vocab_size = len(word_index) + 1
51 embedding_matrix = np.zeros((vocab_size, embedding_dim))
52 for word, i in word_index.items():
53     embedding_vector = embeddings_index.get(word)
54     if embedding_vector is not None:
55         embedding_matrix[i] = embedding_vector
56
57 # Define the Responder model
58 input_layer = Input(shape=(max_sequence_length,))
59 embedding_layer = Embedding(vocab_size, embedding_dim, weights=[embedding_matrix], trainable=
    False)(input_layer)
60 lstm_layer = LSTM(100, return_sequences=True)(embedding_layer)
61 lstm_last_hidden_state = lstm_layer[:, -1, :]
62 lstm_last_cell_state = lstm_layer[:, -1, :]
63 bilstm_layer = Bidirectional(LSTM(100, return_sequences=False), merge_mode='concat')(lstm_layer,
    initial_state=[lstm_last_hidden_state, lstm_last_cell_state, lstm_last_hidden_state,
    lstm_last_cell_state])
64 # bilstm_layer = Bidirectional(LSTM(100, return_sequences=False), merge_mode='concat')(
    embedding_layer, initial_state=[lstm_last_hidden_state, lstm_last_cell_state,
    lstm_last_hidden_state, lstm_last_cell_state])
65 output_layer = Dense(num_classes, activation='softmax')(bilstm_layer)
66 model = Model(inputs=input_layer, outputs=output_layer)
67
68 # model = Sequential()
69 # model.add(Embedding(vocab_size, embedding_dim, input_length=max_sequence_length, weights=[
    embedding_matrix], trainable=False))

```

```

70 # model.add(LSTM(100, return_sequences=True))
71 # model.add(Bidirectional(LSTM(100, return_sequences=True)))
72 # model.add(Bidirectional(LSTM(100)))
73 # model.add(Dense(num_classes, activation='softmax'))
74
75 model.compile(optimizer='adam', loss='sparse_categorical_crossentropy', metrics=['accuracy'])
76
77 # Train the model
78 history = model.fit(X, answers_encoded, epochs=100, batch_size=8)
79
80 # Plot loss and accuracy
81 plt.figure(figsize=(12, 4))
82 plt.subplot(1, 2, 1)
83 plt.plot(history.history['loss'])
84 plt.title('Loss')
85 plt.xlabel('Epochs')
86 plt.ylabel('Loss')
87 plt.subplot(1, 2, 2)
88 plt.plot(history.history['accuracy'])
89 plt.title('Accuracy')
90 plt.xlabel('Epochs')
91 plt.ylabel('Accuracy')
92 plt.tight_layout()
93 plt.savefig('lossaccuracyplot2.pdf')
94 plt.show()
95
96 # Test the model on Table 3 questions
97 test_questions = [
98     'How many people speak French?',
99     'What day is today?',
100     'Who will win the war?',
101     'Who is Italian first minister?',
102     'When World War II ended?',
103     'When Gandhi was assassinated?'
104 ]
105
106 # Normalize test questions
107 test_questions = [normalize_text(question) for question in test_questions]
108
109 test_sequences = tokenizer.texts_to_sequences(test_questions)
110 test_X = pad_sequences(test_sequences, maxlen=max_sequence_length)
111
112 predictions = model.predict(test_X)
113 predicted_labels = [answer_labels[np.argmax(pred)] for pred in predictions]
114

```

```

115 # Save predictions as pandas DataFrame
116 results = pd.DataFrame({'Question': test_questions, 'Predicted Answer': predicted_labels})
117 results.to_csv('predictions.csv', index=False)
118
119 # Print the predicted answers
120 for question, answer in zip(test_questions, predicted_labels):
121     print(f'Question: {question}')
122     print(f'Predicted Answer: {answer}\n')

```

برخی نتایج به صورتی است که در جدول ۶، شکل ۲۵ و برنامه ۴ نشان داده شده است:

Program 4: Prototype Responder: Some Results

```

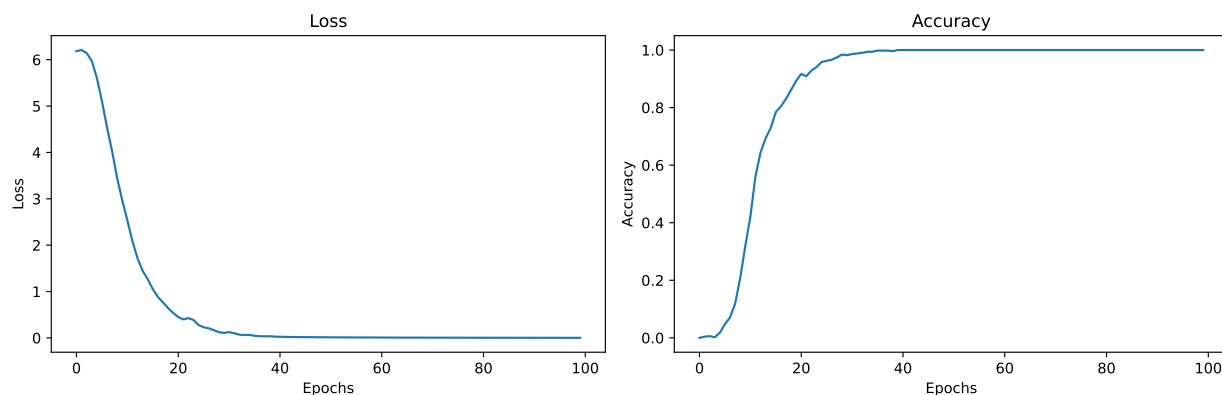
1 Question: how many people speak french?
2 Predicted Answer: 2000
3
4 Question: what day is today?
5 Predicted Answer: tennis tournament
6
7 Question: who will win the war?
8 Predicted Answer: 16th US president
9
10 Question: who is italian first minister?
11 Predicted Answer: Immanuel Kant
12
13 Question: when world war ii ended?
14 Predicted Answer: 1945
15
16 Question: when gandhi was assassinated?
17 Predicted Answer: 1830
18 -----
19 Question: how many people speak french?
20 Predicted Answer: 5,280
21
22 Question: what day is today?
23 Predicted Answer: deliberate killing of a people
24
25 Question: who will win the war?
26 Predicted Answer: Paris
27
28 Question: who is italian first minister?
29 Predicted Answer: 16th US president
30
31 Question: when world war ii ended?
32 Predicted Answer: June 21st
33
34 Question: when gandhi was assassinated?

```


35 Predicted Answer: 1743

جدول ۶: برخی نتایج مربوط به مدل Prototype Responder

Question	Answer
How many people speak French?	2000 5280 ...
What day is today?	tennis tournament
Who will win the war?	16th US president
Who is Italian first minister?	Immanuel Kant
When World War II ended?	1945 June 21st ...
When Gandhi was assassinated?	1830 1743 ...



شکل ۲۵: نمایی از مدل Prototype Responder.