



دانشگاه شهید بهشتی
دانشکده مهندسی و علوم کامپیوتر
پروژه‌های نهایی درس هستان‌شناسی

پروژه اول: پاسخ‌دهی به سوالات دانشکده با رویکرد FAQ و با استفاده از تابع
هزینه زاویه‌ای و پرس و جوی SparQL در هستان‌شناسی

پروژه دوم: پاسخ‌دهی به سوالات دانشکده به کمک یافتن شباهت کسینونی بین
سوالات ورودی و جملات ساخته شده (Predicate, Object, Subject) از هستان -
شناسی

کارشناسی ارشد مهندسی کامپیوتر
گرایش هوش مصنوعی، رباتیکز و رایانش شناختی

دانشجو:
سیدسروش مجد

استاد درس:
خانم دکتر مهنوش شمس‌فرد

تیر ماه ۱۴۰۲

فهرست مطالب

۱- پروژه شماره یک (پاسخدهی به سوالات با رویکرد FAQ و با استفاده از تابع هزینه زاویه‌ای و پرس و جوی SparQL در هستان‌شناسی دانشکده).....	۳
۱-۱ مقدمه	۳
۲-۱ درست کردن مجموعه‌داده سوالات متداول برای آموزش BERT	۳
۳-۱ معرفی تابع هزینه زاویه‌ای	۵
۴-۱ آموزش BERT برای تشخیص کلاس سوالات به کمک تابع هزینه زاویه‌ای.....	۷
۵-۱ تشخیص سوال ورودی، اعمال پیش‌پردازش روی آن و ایجاد پرس و جوی SparQL به ازای هر نوع سوال	۱۲
۶-۱ مزایا و معایب رویکرد	۱۳
۲- پروژه شماره دو پاسخدهی به سوالات دانشکده به کمک یافتن شباهت کسینونی بین سوالات ورودی و جملات ساخته شده (Subject, Object, Predicate) از هستان‌شناسی).....	۱۴
۱-۲ مقدمه	۱۴
۲-۲ درست کردن مجموعه‌داده سوالات برای سنجش عملکرد مدل	۱۵
۳-۲ رویکرد پاسخ به سوالات.....	۱۶
۴-۲ مزایا و معایب	۱۸
۳- مقایسه رویکرد دو پروژه برای پاسخدهی به سوال.....	۱۸

۱- پروژه شماره یک (پاسخ‌دهی به سوالات با رویکرد FAQ و با استفاده از تابع هزینه زاویه‌ای و پرس و جوی SparQL در هستان‌شناسی دانشکده)

۱-۱ مقدمه

در این پروژه به ۲۰ سوال در مورد دانشکده به کمک هستان‌شناسی ساخته شده در گروه پاسخ داده خواهد شد. ابتدا مدل از پیش آموزش دیده BERT انگلیسی را با رویکرد یادگیری متریک و تابع هزینه زاویه‌ای آموزش خواهیم داد. سپس زمان آزمایش ابتدا پیش‌پردازشی بر روی سوالات ورودی اعمال کرده و با استفاده از شبکه BERT آموزش داده شده، حوزه اصلی سوالات ورودی را تشخیص می‌دهیم. در نهایت با توجه به هر سوال پرس و جوی مربوط به آن را در هستان‌شناسی دانشکده انجام داده و جواب آن را به کاربر نشان می‌دهیم.

۱-۲ درست کردن مجموعه داده سوالات متداول برای آموزش BERT

ابتدا یک مجموعه داده‌ای از سوالات متداول ممکن ایجاد و در فایل txt. ذخیره کردیم. این فایل از ۲۰ سوال تشکیل شده که هر سوال به ۱۱ بیان مختلف نوشته شده است. در واقع هر نوع سوال یک کلاس در نظر گرفته شد و به ازای هر سوال با استفاده از ChatGPT ۱۱ نمونه مشابه با نحوه بیان مختلف^۱ ایجاد کردیم. نمونه‌ای از این سوالات را مشاهده می‌کنید:

کلاس سوالات مربوط به درخواست شماره تلفن استاد خاص x:

what is the phone number of dr. x?
can you provide me with dr. x's phone number?
do you know the contact number for dr. x?
how can i reach dr. x by phone?
what is the phone contact of dr. x?
could you give me the phone number for dr. x?

¹ Paraphrase

i need the telephone number of dr. x.

can you share the contact details of dr. x?

what's the phone number for dr. x?

how do i get in touch with dr. x by phone?

how can i call dr. x?

کلاس سوالات مربوط به تعداد واحد مورد نیاز برای مقطع خاص z) masters, bachelors,

:(phd

How many units do I have to pass in z?

how many units are required to pass in z?

what is the unit requirement to complete z?

can you tell me the minimum units needed for z?

what is the total unit count for z?

how many units must I clear in z?

i want to know the units necessary for z.

what is the minimum unit threshold for z?

could you provide information on z unit requirements?

what are the units I need to pass for z?

please inform me about the units required to complete z.

کلاس سوالات مربوط به مکان برگزاری درس خاص c:

which class will the course c be held?

where will course c be held?

at which class is course c scheduled?

what is the class location for course c?

can you provide the class details for course c?

where and which class does course c take place?

what class the course c going to be conducted?

which class is assigned for course c?

what is the class venue for course c?

where can i find information about the class location for course c?

at what class will course c be offered?

کلاس سوالات مربوط به ددلاین دفاع پایان نامه برای مقطع z و ورودی سال t (مثلا مقطع ارشد

ورودی ۱۴۰۰):

when is the deadline for thesis defense for z degree and t graduate level?
what time is the deadline for thesis defense for z degree and t graduate level?
what is the specified deadline for completing the thesis defense for z degree and t graduate level?
can you inform me about the thesis defense deadline for z degree and t graduate level?
when do I need to complete my thesis defense for z degree and t graduate level?
what is the final date to conduct the thesis defense for z degree and t graduate level?
can you provide the deadline for the thesis defense of z degree and t graduate level?
by when should i finish my thesis defense for z degree and t graduate level?
what is the timeframe for the thesis defense requirement of z degree and t graduate level?
can you let me know the latest date for completing the thesis defense for z degree and t graduate level?
what is the cutoff date for the thesis defense for z degree and t graduate level?

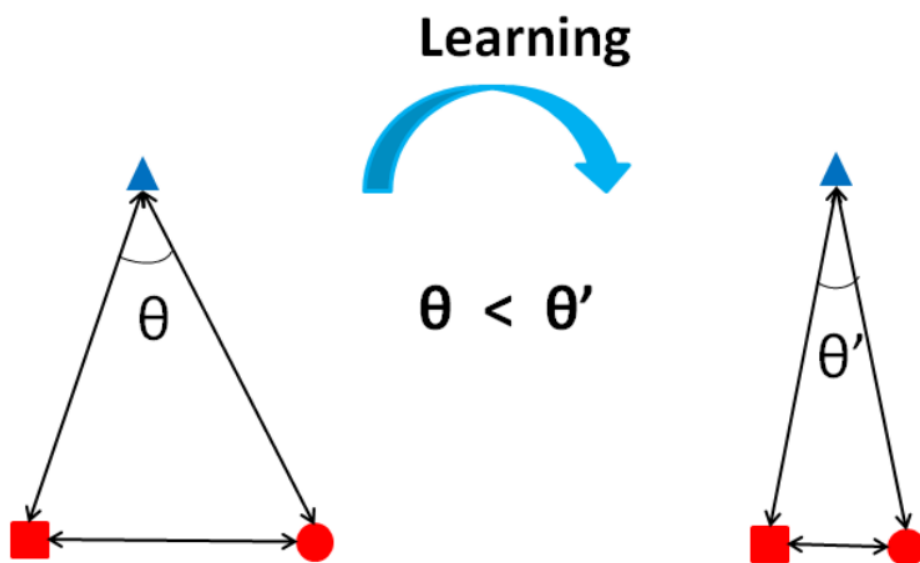
۳-۱ معرفی تابع هزینه زاویه‌ای^۲

تابع هزینه زاویه‌ای پیشنهاد می‌کند بر روی محدودیت زاویه‌ای در نمونه منفی مثلث‌های سه‌گانه (منفی، مثبت و لنگر) تمرکز کنیم. به این صورت که این زاویه از مقدار خاص کمتر باشد. این محدودیت باعث می‌شود فاصله بین مثبت و لنگر کم و فاصله منفی از مثبت و لنگر زیاد شود. در روش‌های قبلی فاصله (لنگر و منفی) و فاصله (مثبت و لنگر) لحاظ می‌شدند ولی در این روش با محدودیت زاویه‌ای هر سه فاصله (لنگر و منفی) و (مثبت و لنگر) و (مثبت و منفی) لحاظ می‌شوند و در نتیجه مقاومت^۳ بیشتر می‌شود (این مورد برای تصویر در مقاله ذکر شده).

² Angular Loss Function

³ Robustness

برای مثال بردارهای جاسازی که با عبور نمونه ورودی از یک شبکه عمیق مانند BERT به دست می‌آیند را در نظر بگیرید. تابع هزینه زاویه‌ای با ایجاد محدودیت برای زاویه تتا، بردارهای جاسازی نمونه‌های متعلق به یک کلاس را به یکدیگر نزدیک و بردارهای جاسازی نمونه‌های کلاس‌های متفاوت را از هم دور می‌کند. نشان داده شده است که تابع هزینه زاویه‌ای در هنگام مواجهه با تغییرات یا اغتشاش‌های زیاد، موثر و مقاوم است. دلیل این مقاومت ایجاد مرز^۴ بیشتر بین نمونه‌های کلاس‌های متفاوت است که زمان آزمایش^۵ با اغتشاش روی نمونه ورودی (مثلاً غلط املائی یا بیان جمله ورودی به شکل دیگر)، جاسازی آن نمونه با جاسازی نمونه‌های کلاس‌های دیگر اشتباه گرفته نشود (یعنی جاسازی آن نمونه به نمونه‌های کلاس خودش نزدیک باشد و به کلاس‌های دیگر نزدیک نباشد).



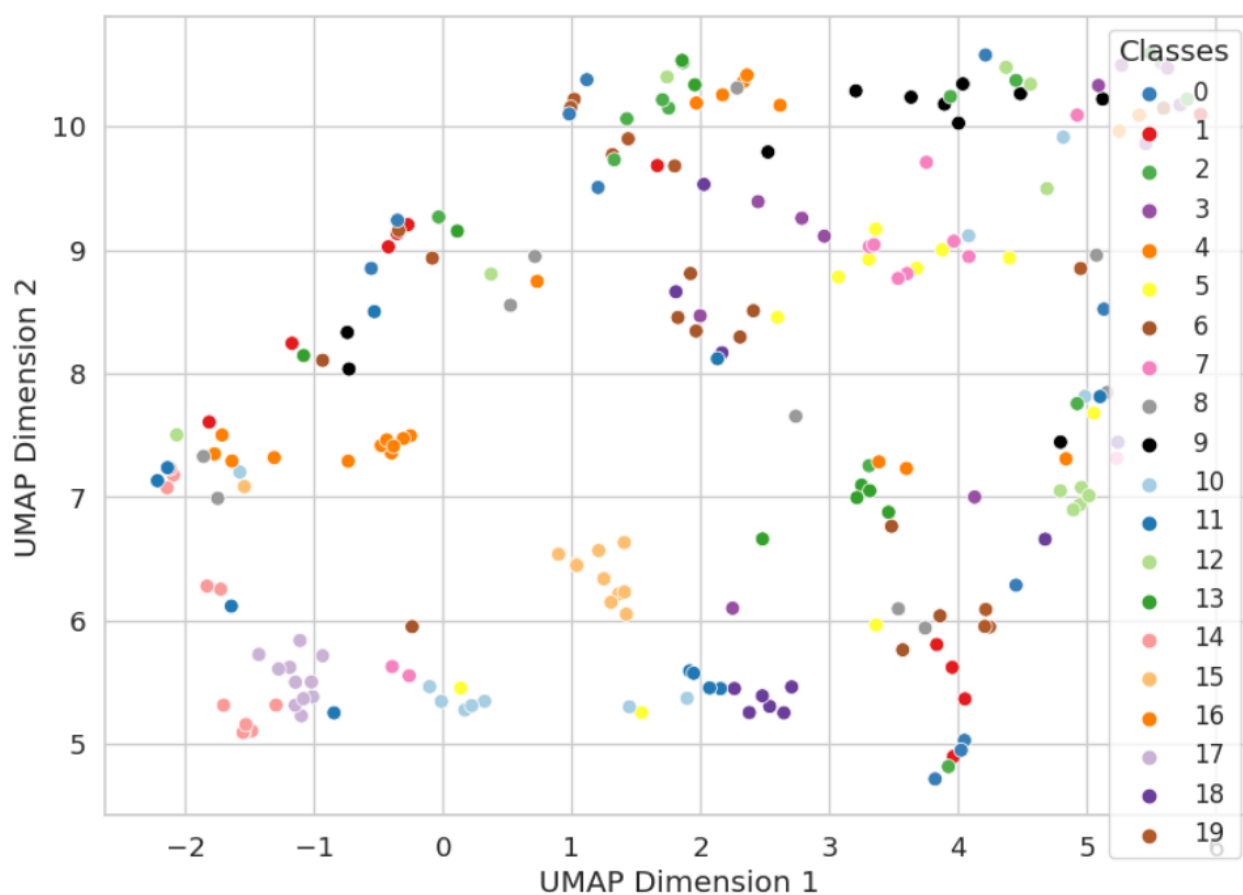
شکل ۱: آموزش شبکه با تابع هزینه زاویه‌ای

^۴ Margin

^۵ Inference

۴-۱ آموزش BERT برای تشخیص کلاس سوالات به کمک تابع هزینه زاویه‌ای

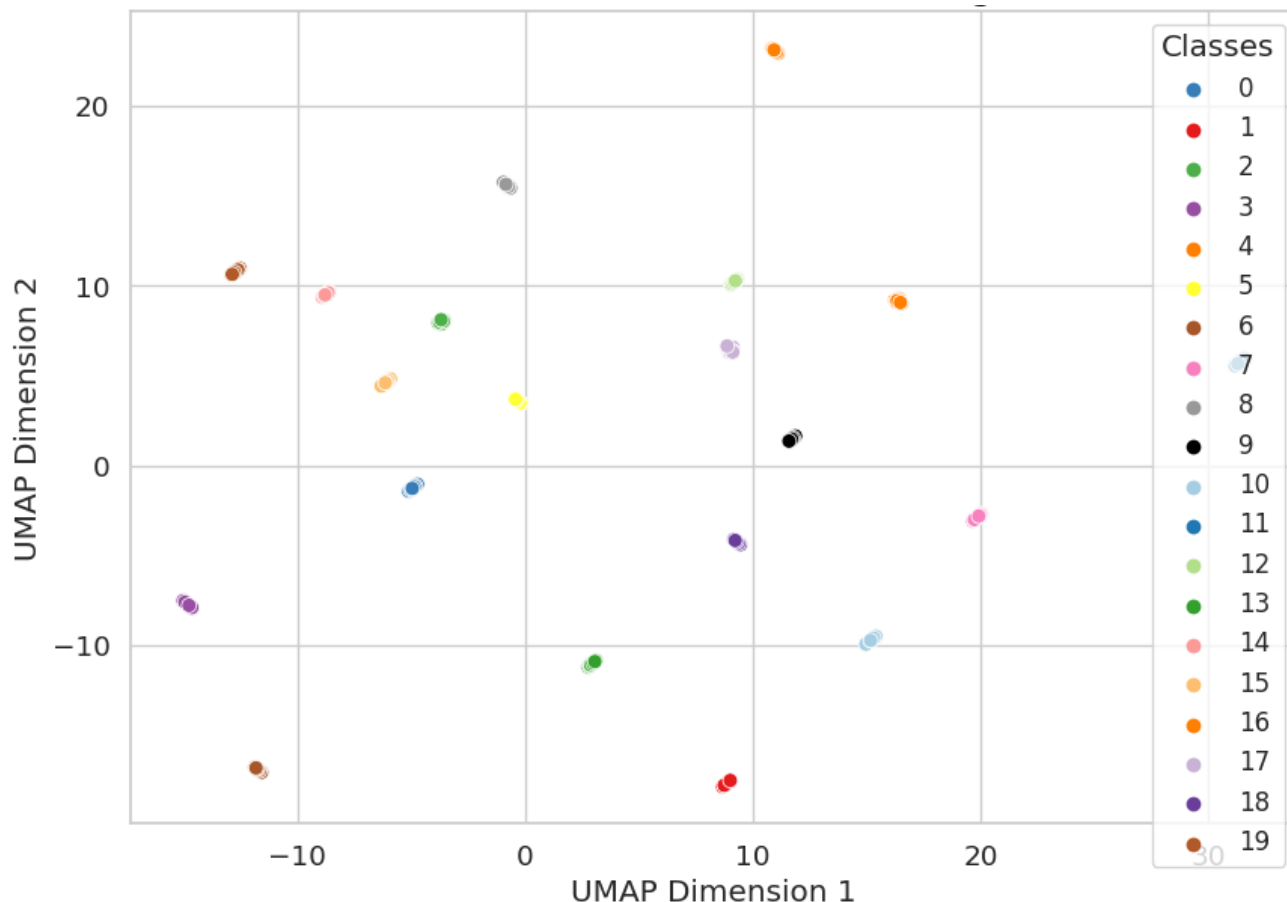
شکل زیر جاسازی BERT از پیش آموزش دیده (هنوز شبکه را با تابع هزینه زاویه ای آموزش ندادیم) بعد از کاهش بعد به کمک Umap را نشان می‌دهد (این روش کاهش ابعاد برای نمایش جاسازی استفاده می‌شود و ساختار زیربنایی و روابط داده‌ها را حفظ می‌کند). هرکدام از نقاط یکی از نمونه‌های سوالات را نمایش می‌دهند و رنگ آن‌ها کلاس آن سوالات را نشان می‌دهد. مثلاً آبی (کلاس ۰) مربوط به سوالات درخواست شماره تلفن استاد خاص است.



شکل ۲: نمایش خروجی شبکه Bert بدون هیچ آموزشی

مشاهده می‌شود که جاسازی خوب و توانایی شناسایی سوالات ورودی را نداریم و در نتیجه شبکه نمی‌تواند تمایز میان سوالات را تشخیص دهد. بعضی از رنگ‌ها مثل نارنجی (نمونه‌های کلاس ۴ سوالات)

نزدیک به هم هستند ولی در کل نمونه‌های کلاس‌های مختلف می‌توانند با یکدیگر اشتباه گرفته شوند. برای آموزش شبکه و شناسایی بهتر سوالات ورودی، یک لایه متراکم^۶ ۲۵۶ تایی به آخر BERT اضافه کرده و آن را با استفاده از تابع هزینه زاویه‌ای آموزش دادیم (بسته^۷ ورودی مدل شامل تمامی ۲۲۰ نمونه آموزشی است و خروجی لایه ۲۵۶ تایی نرمالایز شد) و به جاسازی زیر برای نمونه‌های آموزشی رسیدیم:



شکل ۳: نمایش جاسازی خروجی مدل آموزش دیده با تابع هزینه زاویه‌ای برای داده آموزش

با استفاده از K نزدیک‌ترین همسایه^۸ شبکه را آزمایش^۹ کردیم و به دقت ۱۰۰ درصد برای طبقه‌بندی سوالات در داده آموزش رسیدیم. در نمایش جاسازی به کمک Umap نیز مشاهده می‌شود که سوالاتی

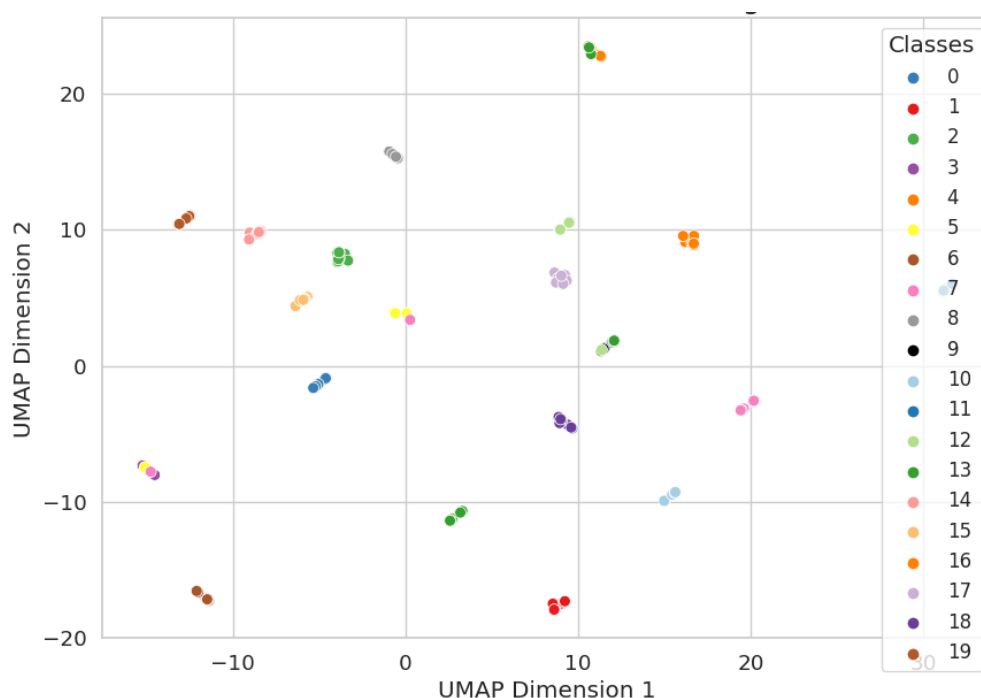
^۶ Dense

^۷ Batch

^۸ K-Nearest-Neighbour

^۹ Test

که مفهوم یکسانی دارند نزدیک به یکدیگر هستند. در واقع شبکه ویژگی‌های متمایز کننده بین کلاس-های مختلف سوالات را یاد گرفته تا این فضای جاسازی خوب به دست آید. نکته حائز اهمیت مقاوم بودن مدل نسبت به بیان مختلف سوالات ورودی و مغشوش کردن آن است. برای مثال اگر غلط املایی داشته باشیم، سوال را به شکل متفاوتی از نمونه‌های آموزشی بیان کنیم، یا متن اضافه‌ای (مثلا احوال پرسی) به اول سوال ورودی اضافه کنیم، اکثر اوقات شبکه خوب عمل کرده و طبقه‌بندی سوال ورودی را به درستی تشخیص می‌دهد. هر ۲۰ سوال متداول را به کمک ChatGPT دوباره به ۱۰ شکل ممکن دیگر نوشتیم (کلا ۲۰۰ نمونه) و مجموعه داده آزمون جدیدی درست کردیم که دارای غلط املایی نیز هستند. همچنین به بعضی از سوالات متونی از احوال پرسی قبل از پرسش سوال و تشکر بعد از پرسش سوال اضافه شد (مانند hello. may i have the contact detel of dr. x, like the phone numbr? many thanks). در نهایت با umap جاسازی خروجی مدل به ازای نمونه‌های آزمون به صورت شکل زیر به دست آمد:



شکل ۴: نمایش جاسازی خروجی مدل آموزش دیده با تابع هزینه زاویه‌ای برای داده‌های آزمون

Accuracy برای مجموعه داده داده آزمون برابر با ۹۶٪ و برای داده آموزش ۱۰۰٪ محاسبه شد. نمونه-

های زیر را مدل به درستی شناسایی نکرده است:

Incorrect prediction:

Predicted label: 3

Ground truth label: 5

Sample: what is the limit on the number of semesters for the z degree?

Incorrect prediction:

Predicted label: 3

Ground truth label: 5

Sample: could you provide information about the z degree's maximum semester count?

Incorrect prediction:

Predicted label: 3

Ground truth label: 7

Sample: can you tell me the minimum required essential courses to get a z degree?

Incorrect prediction:

Predicted label: 5

Ground truth label: 7

Sample: hello, how's your mood today? how many courses are considered essential to complete a z degree?

Incorrect prediction:

Predicted label: 9

Ground truth label: 12

Sample: can you direct me to the place where I can find book b?

Incorrect prediction:

Predicted label: 4

Ground truth label: 13

Sample: hi, how's your state of mind? I'm looking for the specific place where class c1 is held. I'm truly touched by your kindness.

Incorrect prediction:

Predicted label: 9

Ground truth label: 13

Sample: where is the designated room for class c1?

Incorrect prediction:

Predicted label: 4

Ground truth label: 13

Sample: hi, how's your day shaping up? hey, how are you coping? could you direct me to the place where class c1 is conducted?

مشاهده می‌شود وقتی کلمات مهم و حساس که در پیشبینی مدل تاثیر زیادی دارند را با شکل املایی غلط می‌نویسیم مدل اشتباه می‌کند. مثلاً وقتی تعداد ترم مجاز برای مقطع ارشد مورد سوال است و به جای semester، semstr نوشته شده است، مدل پیشبینی غلط کرده. یا وقتی آدرس کلاس خاصی سوال شده است و به جای class، clas نوشتیم مدل اشتباه کرده. البته برای همه‌ی کلمات مهم این اشتباه رخ نداده است و همچنین مدل نسبت به متون اضافه (مانند احوال پرسی و تشکر) مقاوم می‌باشد و این جملات در پیشبینی مدل تاثیرگذار نیستند.

دقت طبقه‌بندی شبکه بسیار به داده آموزشی مرتبط است و اگر داده آموزشی خوبی وجود داشته باشد، عملکرد شبکه بهتر می‌شود. برای این کار می‌توان با ChatGPT نمونه‌های متنوع‌تر و بیشتری برای هر کلاس سوال تولید کرد و به راحتی مدل را با آن داده‌ها مجدداً آموزش داد.

۵-۱ تشخیص سوال ورودی، اعمال پیش‌پردازش روی آن و ایجاد پرس و جوی SparQL به ازای هر نوع سوال

برای پاسخ دهی به سوال ورودی ابتدا پیش‌پردازش اولیه‌ای بر روی آن جهت استخراج مواردی مانند اسامی اساتید، نام دانشکده، نام دروس، نام کتاب‌ها، شماره کلاس و مقطع تحصیلی انجام دادیم و به جای آن‌ها به ترتیب مقادیر x, y, c, b, cl و z را جایگذاری کردیم. این کار را برای این انجام می‌دهیم که بتوانیم حوزه اصلی سوال را با شبکه BERT آموزش دیده با تابع هزینه زاویه‌ای تشخیص دهیم و بفهمیم سوال از ما چه می‌خواهد. برای مثال جمله زیر را در نظر بگیرید:

Dr. shamsfard is teaching nlp for bachelors degree in the computer engineering department. The book russel is recommended for class 101.

جمله به صورت زیر تغییر داده خواهد شد (این جمله صرفاً برای نمایش عملکرد پیش پردازش می‌باشد). و آن را به شبکه BERT می‌دهیم:

: Dr. x is teaching c for z degree in the y department. The book b is recommended for class cl.

و دیتافریمی به شکل زیر به دست خواهد آمد:

x	y	c	b	cl	z
shamsfard	computer engineering	nlp	russel	101	bachelors

عبارت سوال که تغییر داده شده را به مدل آموزش دیده می‌دهیم و بردار جاسازی آن را به دست می‌آوریم. سپس با Transform کردن K نزدیک‌ترین همسایه (که بعد از آموزش مدل روی داده‌های آموزش Fit کرده بودیم) بر بردار جاسازی عبارت ورودی، حوزه اصلی سوال ورودی را تشخیص می‌دهیم (منظور همان کار طبقه‌بندی سوال از ۲۰ سوال ممکن است). بعد از آن که حوزه اصلی سوال تشخیص داده شد، با استفاده از پرس و جوی SparQL و موجودیت‌های نام دار موجود در دیتافریم، پاسخ آن را در هستان‌شناسی می‌یابیم. در واقع به ازای هر سوال، به ازای مقادیر مختلفی که دیتافریم می‌تواند داشته باشد روال پرس‌وجوی SparQL نوشته شده است. با توجه به اینکه هستان‌شناسی تولید شده توسط گروه به زبان فارسی و سیستم تشخیص سوالات به زبان انگلیسی است، مجبور شدم در کد تعداد زیادی if بنویسم. در صورتی که زبان هستان‌شناسی و سیستم پاسخ‌دهی به سوالات یکسان باشد یا هستان‌شناسی چند زبانه باشد، روال پاسخ‌دهی به سوالات بیشتر خودکار خواهد شد و نیازی به نوشتن if زیاد نخواهد بود. به این صورت که اسامی Instance‌های مفاهیم مختلف را در پرس و جوی SparQL می‌نویسیم. برای مثال در بخش پیش پردازش به جای اسامی Instance‌های کلاس استاد در جملات، x جایگذاری

کرده و در پرسوجوی SparQL نیز نام آن Instance ای که به جای آن x جایگذاری شده است را جست و جو می‌کنیم. در صورتی که هستان‌شناسی زبان انگیزی نیز داشته باشد با تغییر کمی در کد می‌توان خودکارسازی بیشتری انجام داد.

۱-۶ مزایا و معایب رویکرد

• مزایا

- ۱- در صورتی که بخواهیم به سوالات بیشتری پاسخ دهیم می‌توان نمونه آن سوالات را به مجموعه داده اضافه کرده و مدل را بدون تغییر معماری آموزش داد.
- ۲- متون اضافه در سوالات مانند تشکر و احوال‌پرسی اولیه بر رویکرد ارائه شده تاثیر منفی ندارد (با توجه به اینکه در مجموعه داده آموزش این متون اضافه وجود ندارد). به بیان دیگر مدل اصل مطلب و سوال پرسیده شده را استخراج می‌کند.
- ۳- مدل مقاومت خوبی نسبت به غلط‌های املائی نشان می‌دهد.
- ۴- روال پاسخ‌دهی به هر سوال مشخص است و بعد از تشخیص درست سوال قطعا جواب آن درست داده می‌شود.

• معایب

- ۱- در رویکرد FAQ باید از قبل نوع سوالات و روال پاسخ‌دهی به آن‌ها مشخص باشد و از این لحاظ محدودیت داریم.
- ۲- با توجه به اینکه زبان هستان‌شناسی فارسی و تشخیص سوال انگیزی است، مجبور شدیم if‌های زیادی استفاده کنیم و اگر به هستان‌شناسی instance‌های بیشتری اضافه شود، مدل توانایی عملکرد خوب برای instance‌های اضافه شده را نخواهد داشت. اگر هستان‌شناسی انگیزی باشد این مشکل با تغییر کمی در کد رفع می‌شود.
- ۳- مدل برای برخی سوالات با کلمات خاص دارای غلط املائی درست پیشبینی نمی‌کند (در صورتی که جاسازی در سطح کاراکتر به مدل اضافه شود می‌توان این مشکل را رفع کرد یا برای مثال با استفاده از جاسازی کلمات مهم و حساس، قبل از ورود جمله به شبکه غلط‌های املائی را اصلاح و با کلمات درست کرد).

(مجموعه داده آموزش و آزمایش در فایل زیپ ضمیمه شده است)

۲- پروژه دوم (پاسخ‌دهی به سوالات دانشکده به کمک یافتن شباهت کسینونی بین سوالات ورودی و جملات ساخته شده (Object, Subject, Predicate) از هستان‌شناسی)

۲-۱ مقدمه

در این پروژه می‌خواهیم با استفاده از جاسازی هستان‌شناسی دانشکده، به سوالات پاسخ دهیم. جملات را به صورت سه تایی subject ، object و predicate که از هستان‌شناسی استخراج شده اند نوشته و جاسازی این جملات را با parsbert به دست می‌آوریم (پیش‌پردازش اولیه روی آن‌ها انجام گرفت مثلاً " _"ها حذف شدند). سپس با یافتن جوابی که بیشترین شباهت به سوال ورودی دارد به سوالات پاسخ می‌دهیم. برای سنجش عملکرد مدل ، ۲۱ سوال طراحی شد با استفاده از ChatGPT هر کدام از سوالات به ۱۰ بیان مختلف نوشته شد.

۲-۲ درست کردن مجموعه داده سوالات برای سنجش عملکرد مدل

برای سنجش عملکرد مدل ، ۲۱ سوال طراحی شد با استفاده از ChatGPT هر کدام از سوالات به ۱۰ بیان مختلف نوشته شد. چند نمونه از سوالات را در زیر مشاهده می‌کنید:

سوال:

- "دکتر شمس فرد از کجا فارغ التحصیل شده است؟"
- "از کدام دانشگاه دکتر شمس فرد فارغ التحصیل شده است؟"
- "دکتر شمس فرد در کدام مرکز تحصیلی فارغ التحصیل شده است؟"
- "دکتر شمس فرد از چه موسسه‌ای فارغ التحصیل شده است؟"
- "محل فارغ التحصیلی دکتر شمس فرد کجاست،"
- "از کدام موسسه تحصیلی دکتر شمس فرد فارغ التحصیل شده است؟"
- "دکتر شمس فرد در کدام دانشگاه فارغ التحصیل شده است؟"
- "محل فارغ التحصیل شدن دکتر شمس فرد کجا بوده است"
- "دکتر شمس فرد از کجا فارغ التحصیل شده است؟"
- "از کدام مرکز تحصیلی دکتر شمس فرد فارغ التحصیل شده است؟"

جواب:

دکتر شمس فرد محل فارغ التحصیلی دانشگاه صنعتی امیرکبیر.

سوال:

- "مقدار مبلغ شهریه ثابت در مقطع ارشد چند ریال است؟"
- "مبلغ ثابت شهریه در دوره ارشد چقدر است؟"
- "مبلغ شهریه ثابت در مقطع ارشد چند ریال است؟"
- "چند ریال برای مبلغ شهریه ثابت در دوره ارشد تعیین شده است؟"
- "مقدار مبلغ شهریه ثابت در مقطع ارشد چقدر می باشد؟"
- "چه مبلغی برای شهریه ثابت در مقطع ارشد تعیین شده است؟"
- "مبلغ ثابت شهریه در مقطع ارشد چه مقداری است؟"
- "مبلغ شهریه ثابت در دوره ارشد چقدر می باشد؟"
- "چه مبلغی برای مبلغ شهریه ثابت در مقطع ارشد مشخص شده است؟"
- "مقدار مبلغ شهریه ثابت در دوره ارشد چند ریال می باشد؟"

جواب:

قانون شهریه ثابت ارشد مبلغ شهریه ۵۹۰۰۰۰۰۰ ریال.

سوال:

- "ساعت حضور دکتر قوامی زاده چه موقع است؟"
- "دکتر قوامی زاده در چه ساعاتی حاضر است، پ"
- "ساعت حضور دکتر قوامی زاده در چه زمانی است؟"
- "چه ساعت‌هایی دکتر قوامی زاده حاضر است؟"
- "ساعت حضور دکتر قوامی زاده در کدام بازه زمانی است؟"
- "دکتر قوامی زاده در چه بازه‌های زمانی حاضر است؟"
- "چه زمانی برای ساعت حضور دکتر قوامی زاده مشخص شده است؟"
- "دکتر قوامی زاده در کدام ساعات حاضر است؟"
- "ساعت حضور دکتر قوامی زاده در چه بازه‌های زمانی است؟"
- "چه ساعاتی دکتر قوامی زاده حاضر است؟"

جواب:

دفتر دکتر قوامی زاده ساعت حضور شنبه تا سه شنبه ساعت هشت تا نه.

۲-۳ رویکرد پاسخ به سوالات

ابتدا تمامی جملات ممکن را با تمامی subject, object و predicate موجود در هستان‌شناسی می‌سازیم. و جاسازی همه آن‌ها را با مدل ParsBERT به دست می‌آوریم (در این رویکرد شبکه را آموزش نمی‌دهیم). برای پاسخ به سوال ورودی، جاسازی عبارت سوال ورودی را با همان مدل ParsBERT به دست آورده و فواصل کسینوسی این جاسازی را با تمامی جاسازی‌های جملات ساخته شده محاسبه می‌کنیم. سپس جمله با بیشترین شباهت به سوال ورودی (کمترین فاصله کسینوسی) را به عنوان پاسخ نهایی بر می‌گردانیم. برای سنجش عملکرد مدل، معیار Accuracy را به صورت Top-1، Top-3 و Top-5 نمایش می‌دهیم (در معیار top-3 در صورتی که فاصله کسینوسی جواب مورد نظر در بین سه تا از کمترین فواصل کسینوسی با سوال ورودی باشد، یعنی شبکه درست پیشبینی کرده. به همین منوال برای top-5 اگر بین پنج تا از کمترین فواصل باشد یعنی شبکه به درستی پیشبینی کرده است). در نهایت برای همه سوالات (۲۰۰ نمونه) Accuracy به صورت زیر محاسبه شد:

جدول ۱: معیارهای Accuracy برای سنجش عملکرد مدل در پاسخ‌دهی به سوالات

Top-1 Accuracy	Top-3 Accuracy	Top-5 Accuracy
۴۸.۰۹٪	۷۶.۶۶٪	۸۵.۷۱٪

بیشترین دقت:

در ارشد چند ترم می‌توان مرخصی گرفت؟
تعداد ترم‌های مجاز برای مرخصی در دوره ارشد چند است؟
در دوره ارشد چه تعداد ترم مرخصی قابل استفاده است؟
در ارشد به چه تعداد ترم مرخصی اجازه داده می‌شود؟
تعداد ترم‌های قابل استفاده برای مرخصی در دوره ارشد چه مقداری است؟
در دوره ارشد چند ترم مرخصی قابل استفاده است؟
چه تعداد ترم برای مرخصی در دوره ارشد مجاز است؟
در دوره ارشد به چه تعداد ترم مرخصی اجازه داده می‌شود؟
تعداد ترم‌های مجاز برای مرخصی در ارشد چند است؟
در ارشد چند ترم مرخصی قابل استفاده است؟

Top-1 Accuracy: 100.00%

Top-3 Accuracy: 100.00%

Top-5 Accuracy: 100.00%

پاسخ: قانون مرخصی گرفتن ارشد تعداد ترم ۱.

کمترین دقت:

ساعت حضور دکتر شمس فرد چه موقع است؟
دکتر شمس فرد در چه ساعاتی حاضر است؟
ساعت حضور دکتر شمس فرد در چه زمانی است؟
چه ساعت‌هایی دکتر شمس فرد حاضر است؟
ساعت حضور دکتر شمس فرد در کدام بازه زمانی است؟
دکتر شمس فرد در چه بازه‌های زمانی حاضر است؟
چه زمانی برای ساعت حضور دکتر شمس فرد مشخص شده است؟
دکتر شمس فرد در کدام ساعات حاضر است؟
ساعت حضور دکتر شمس فرد در چه بازه‌های زمانی است؟
چه ساعاتی دکتر شمس فرد حاضر است؟

Top-1 Accuracy: 0.00%
Top-3 Accuracy: 10.00%
Top-5 Accuracy: 10.00%

پاسخ: دفتر دکتر شمس فرد ساعت حضور شنبه‌ها از ساعت ده تا دوازده، چهارشنبه‌ها از ساعت سه تا پنج.

تمامی سوالات مجموعه داده و دقت مدل برای هر کدام در فایل نوت‌بوک موجود و قابل مشاهده

است.

۲-۴ مزایا و معایب

• مزایا

- ۱- در این رویکرد بدون نیاز به آموزش مدل می‌توانیم به سوالات ورودی پاسخ دهیم.
- ۲- برای هستان‌شناسی‌های مختلف (حتی حوزه‌های مختلف) می‌توانیم این رویکرد را استفاده کنیم و نیازی به کار اضافه نیست. در نتیجه نسبت به رویکرد FAQ محدودیت کمتری دارد.

• معایب

- ۱- در صورتی که هستان‌شناسی دارای روابط اضافه باشد، ممکن است مدل به درستی کار نکند. برای مثال در یکی از سوالات تعریف قانون آموزش محور خواسته شده است و جوابش به صورت زیر است:
"قانون آموزش محور برای ارشد تعریف آموزش محور اگر تا آخر ترم سه، هجده واحد گذرانده نشود، دانشجوی کارشناسی ارشد، آموزش محور میشود."

یکی از جملاتی که از سه‌تایی‌های RDF هستان‌شناسی دانشکده ساخته شده "قانون آموزش محور برای ارشد وضع میشود برای دانشگاه شهیدبهشتی." است. برای پاسخ به سوال تعریف قانون آموزش محور شبکه این جمله سبز را به عنوان جواب به اشتباه بر می‌گرداند. بنابراین یکی از معایب مهم این رویکرد این است که باید در ساخت هستان‌شناسی توجه داشته باشیم روابط را طوری تعریف کنیم که وقتی

جملات از سه تایی های RDF ساخته شدند، شبیه به یکدیگر نباشند تا در پاسخ دهی به سوالات، مدل اشتباه نکند. این اشتباه در این رویکرد بسیار پررخداد بود.

۲- با توجه به اینکه مدل هیچ آموزشی ندیده، نسبت به غلط های املایی و نحوه مختلف بیان سوالات مقاوم نیست .

۳- مقایسه رویکرد دو پروژه برای پاسخ دهی به سوال

در پروژه اول نسبت به اغتشاش در سوالات ورودی مقاومت خوبی ایجاد شد ولی سوالات و روال پاسخ به آن ها باید از قبل مشخص شده باشند. همچنین شبکه باید آموزش داده شود. در پروژه دوم نیازی به این کارها نبود و بدون تعریف سوالات از قبل و بدون نیاز به آموزش مدل می توانستیم به سوالات مربوط به هستان شناسی پاسخ دهیم. ولی نسبت به رویکرد پروژه اول مدل به اغتشاش در سوال ورودی حساس تر و دقت شناسایی سوال آن نیز کمتر است. در شکل ۲ برای جاسازی مدل BERT دیدیم که BERT خیلی عملکرد خوبی بدون اینکه آموزش ببیند ندارد.

با توجه به کاربرد، نیاز، نوع هستان شناسی، انرژی و زمانی که در دسترس هست می توان از هر دو رویکرد استفاده کرد. برای مثال اگر مقاومت نسبت به اغتشاش در سوال ورودی مهم است رویکرد پروژه اول گزینه مناسب تری است و اگر تعمیم پذیر بودن به هستان شناسی های مختلف اهمیت داشته باشد رویکرد پروژه دوم بهتر است. همچنین در پروژه دوم می توان بعد از ساختن جملات با سه تایی های RDF ، آن ها را Augment کرد و با استفاده از توابع هزینه یادگیری متریک جاسازی بهتری برای ParsBERT به دست آورد (این رویکرد پیاده سازی نشد ولی به نظر می رسد موثر باشد).