

# بسم الله الرحمن الرحيم



تمرین سری سوم

یافتن عبارات مشابه در متون دینی با استفاده از جاگذاری کلمات

پردازش زبان‌های طبیعی

بهار ۱۴۰۱

---

۳	..... شناسنامه
۴	..... تعریف مسئله
۵	..... داده‌ها و پیش‌پردازش
۶	..... روش اول – یافتن عبارات مشابه با استفاده از TF-IDF
۱۳	..... روش دوم – استفاده از Fasttext
۲۰	..... روش سوم – استفاده از شبکه‌های Fine-tuned BERT
۲۰	..... آموزش دادن مدل برت



# یافتن عبارات مشابه در متون دینی با استفاده از جاگذاری کلمات

تکلیف سری سوم

شناسنامه

تمرین سری سوم	
دانشگاه صنعتی شریف – دانشکده کامپیوتر	
۱۴۰۰-۱۴۰۱	
مشخصات درس	
عنوان درس	پردازش زبان‌های طبیعی
استاد مربوطه	آقای دکتر احسان‌الدین عسگری
استاد حل تمرین	تیم اساتید حل تمرین
مشخصات گروه	
نام و نام خانوادگی	امیر پورمند پویا خانی مهدی آخی
مشخصات نویسنده و ارسال کننده گزارش	مهدی آخی – ۹۹۲۰۱۴۹۸
نشانی الکترونیکی	mahdiakhi@ce.sharif.edu
مشخصات سند	
عنوان	گزارش تمرین سری سوم
تاریخ تحویل	۱۴۰۱/۳/۱۹

## تعریف مسئله

برای این که بتوان از علوم داده در پردازش متن و به طور خصوص پردازش متون زبان طبیعی استفاده کرد نیاز است تا آنها را به فرمی قابل استفاده و قابل درک برای کامپیوترها تبدیل کرد نیاز است تا حروف و علائم نگارشی به عدد تبدیل شوند (Word to Vec). برای این کار روش‌های مختلفی وجود دارد که به صورت کلی به دودسته تقسیم می‌شوند: دسته اول روش‌هایی هستند که خاصیت‌های آماری کلمات را به عدد مدل می‌کنند مانند TF-IDF و دسته دوم روش‌هایی هستند که خواص معنایی کلمات را در خود دارند مانند FastText یا BERT. هر دوی این دسته‌ها کاربردهای گوناگونی دارند مانند حدس کلمه، تکمیل جمله، ترمیم جملات، پیدا کردن جملات مشابه (متنی و معنایی). در این تکلیف سعی داریم تا با بهره‌گیری از هر دو دسته عملیات پیدا کردن عبارات مشابه در متون دینی (قرآن کریم، نهج البلاغه و صحیفه سجادیه) را پیاده‌سازی کنیم. آن چه در ادامه این گزارش آمده شرح عملیات و تکنیک‌های استفاده شده در این مسیر است.

# یافتن عبارات مشابه در متون دینی با استفاده از جاگذاری کلمات

تکلیف سری سوم

## داده‌ها و پیش‌پردازش

تمامی داده‌ها از گیت‌هاب درس به آدرس زیر برداشته شده‌اند.

```
git clone https://github.com/language-ml/course-nlp-ir-1-text-exploring
```

لود کردن آیات قرآن و عبارات نهج البلاغه چالش خاصی نداشت ولی لود کردن عبارات صحیفه سجادیه نیاز به پیش‌پردازش‌های دیگری هم داشت زیرا آیات با شماره آیه شروع می‌شدند و بعضی از خطوط نیز اصلاً شامل عبارتی نبودند. سپس داده‌ها با استفاده از کتابخانه camel\_tools که یک کتابخانه استاندارد برای پردازش متن در زبان عربی بصورت عام هست، پیش‌پردازش شدند. از آنجایی که متون ثابتی در مجموعه داده وجود داشت عملیات پیش‌پردازش شامل مراحل عمومی بود. از جمله:

- حذف اعراب و تنوین‌ها
- حذف برخی علائم قرائت
- حذف کلمات با اندازه کوچکتر از ۳
- یک شکل کردن حروفی مانند «ک» و «ی» عربی

. در اینجا یک نمونه از داده‌های پیش‌پردازش شده برای هر کدام از ۳ رفرنس اصلی نشان داده شده است.

```
*Quranic(pure): وَإِذَا سَأَلَكَ عِبَادِي عَنِّي فَإِنِّي قَرِيبٌ ۚ أُجِيبُ دَعْوَةَ الدَّاعِ إِذَا دَعَانِ ۚ فَلْيَسْتَجِيبُوا لِي وَلْيُؤْمِنُوا بِي لَعَلَّهُمْ يَرْشُدُونَ
*Quranic(Processed): واذا سالك عبادي عني فاتي قريب اجيب دعوه الداع اذا دعان فليستجيبوا لي وليؤمنوا بي لعلهم يرشدون

*Nahj(pure): وَقَالَ عَ لِلظَّالِمِ الْبَادِي غَدًا بِكَفَّةٍ عَصَاهُ
*Nahj(Processed): و قال ع للظالم البادي غدا بكفه عصه

*Sahifa(pure): حَمْدًا يَرْتَفِعُ مِنَّا إِلَىٰ أُطَىٰ جَبَّتَيْنِ فِي كِتَابٍ مَرْفُومٍ يَنْهَدُهُ الْمُفَرِّقُونَ .
*Sahifa(Processed): و كاشف في الدعاء اليك حاتمته
```

بعد از عملیات پیش‌پردازش نوبت به توکن کردن جملات (داکیومنت‌ها) می‌رسد. در تصویر زیر خروجی نمونه این مرحله را مشاهده می‌کنید. می‌توانید ببینید که در بردار زیر کلمات کمتر از سه حرف حذف شده‌اند.

```
Tokenized:
*Quranic Example: ['قَاتِي', 'الْم', 'رَبِّ', 'كَلِمَات', 'قَلْب', 'طِي', 'اِنَّ', 'الشَّوَاب', 'الرَّحِيم']
*Nahj Example: ['فَج', 'الله', 'مُصَلِّه', 'فَل', 'فَل', 'السَّه', 'فَرَار', 'العِيد', 'فما', 'الطَّق', 'مَلَح', 'خَي', 'اسْكَن', 'صَلَّى', 'وَصَفَه', 'خَي', 'بَكَن', 'اَقَام', 'لَاخَنَّا', 'مِسْرَه', 'اَنْظَرْنَا', 'بَمَلَه', 'وَفَرَه']
*Sahifa Example: ['فَج', 'الله', 'مُصَلِّه', 'فَل', 'فَل', 'السَّه', 'فَرَار', 'العِيد', 'فما', 'الطَّق', 'مَلَح', 'خَي', 'اسْكَن', 'صَلَّى', 'وَصَفَه', 'خَي', 'بَكَن', 'اَقَام', 'لَاخَنَّا', 'مِسْرَه', 'اَنْظَرْنَا', 'بَمَلَه', 'وَفَرَه']
```

## روش اول – یافتن عبارات مشابه با استفاده از TF-IDF

در این بخش در ابتدا از داده‌های پیش‌پردازش شده در مرحله قبل استفاده می‌کنیم به طوری که آنها را توسط تابع نرمال‌سازی که داشتیم، نرمال می‌کنیم.

```
#Normalize all verses from quran, Nahj, and Sahife
#
verse_dict_nrmlz = {k:normalize_arabic(v) for k,v in tqdm.tqdm(verse_complete_dict.items())}
nahj_dict_nrmlz = {k:normalize_arabic(v) for k,v in tqdm.tqdm(nahj_complete_dict.items())}
sahife_dict_nrmlz = {normalize_arabic(v) for v in tqdm.tqdm(sahife_complete_dict)}
#
100%|██████████| 6236/6236 [00:00<00:00, 11174.55it/s]
100%|██████████| 800/800 [00:00<00:00, 3515.74it/s]
100%|██████████| 924/924 [00:00<00:00, 7923.86it/s]
```

سپس دیکشنری‌های نرمال‌شده و نرمال‌نشده را به لیست تبدیل کرده و سپس لیست‌ها را به دیتافریم‌های کتابخانه Pandas تبدیل می‌کنیم تا بتوانیم ساده‌تر با آنها کار کنیم.

```
# Converting standard and normalize dictionaries(index:verse) from quran, Nahj, and Sahife to Pandas DataFrame
# with column names: num and verse
#
df_verse_complete_dict = pd.DataFrame(list(verse_complete_dict.items()), columns = ['num','verse'])
df_verse_complete_dict_nrmlz = pd.DataFrame(list(verse_dict_nrmlz.items()), columns = ['num','verse'])
df_nahj_complete_dict = pd.DataFrame(list(nahj_complete_dict.items()), columns = ['num','verse'])
df_nahj_complete_dict_nrmlz = pd.DataFrame(list(nahj_dict_nrmlz.items()), columns = ['num','verse'])
df_sahife_complete_dict = pd.DataFrame(list(sahife_complete_dict), columns = ['verse'])
df_sahife_complete_dict_nrmlz = pd.DataFrame(list(sahife_dict_nrmlz), columns= ['verse'])
#
```

سپس توسط کتابخانه scikit-learn به تعداد ۶ تا آبجکت از وکتورایزهای TF-IDF می‌سازیم و توسط نام‌گذاری آن‌ها را برای استفاده‌های بعدی متمایز می‌کنیم.

```
# Create TfidfVectorizer object (first three ones are Word level and others are Character Level)
#
vectorizer_quran = TfidfVectorizer(lowercase = False)
vectorizer_nahj = TfidfVectorizer(lowercase = False)
vectorizer_sahife = TfidfVectorizer(lowercase = False)
vectorizer_quran_char = TfidfVectorizer(lowercase = False,analyzer='char')
vectorizer_nahj_char = TfidfVectorizer(lowercase = False,analyzer='char')
vectorizer_sahife_char = TfidfVectorizer(lowercase = False,analyzer='char')
#
```

## یافتن عبارات مشابه در متون دینی با استفاده از جاگذاری کلمات

تکلیف سری سوم

سپس این آبجکت‌ها را روی ستون verse که ستونی هست که آیه‌ها و عبارات کتب دینی را در خود دارند، به عبارتی fit کرده و سپس داده‌ها را به آن embedding تبدیل می‌کنیم. خروجی این کار ۶ ماتریس می‌شود که با توجه به ویژگی‌های مشخص شده، ابعاد متفاوتی دارند. این ماتریس‌ها به تعداد آیه/عبارت کتب دینی سطر و به تعداد واحدهای متفاوت و متمایز در این متون ستون دارند. برای مثال ۳ ماتریس اول مربوط به بدست آوردن تعبیه برای کتب قرآن، نهج‌البلاغه و صحیفه سجادیه است به طوری که روی کلمات این کتب این فضای تعبیه ساخته شود. لذا برای قرآن به تعداد کلمات منحصر به فرد آن یعنی ۱۴۶۵۸ تا، برای نهج‌البلاغه نیز به همین منوال ۲۱۰۹۴ تا و برای صحیفه سجادیه نیز ۷۷۱۱ ستون در ماتریس متناظرشان در نظر گرفته می‌شود. ۳ ماتریس بعدی متعلق به همین ۳ کتاب اما در حالتی هستند که روی کاراکترهای متون تعبیه ساخته می‌شود. لذا برای قرآن ۴۰، نهج‌البلاغه ۳۶ و صحیفه سجادیه ۴۷ ستون در نظر گرفته می‌شود (بخاطر کاراکترهای خاص در متون تعداد ستون‌ها متفاوت است).

```
# Generate matrices of word vectors
#
tfidf_matrix_quran = vectorizer_quran.fit_transform(df_verse_complete_dict_nrmzl['verse'])
tfidf_matrix_nahj = vectorizer_nahj.fit_transform(df_nahj_complete_dict_nrmzl['verse'])
tfidf_matrix_sahife = vectorizer_sahife.fit_transform(df_sahife_complete_dict_nrmzl['verse'])
tfidf_matrix_quran_char = vectorizer_quran_char.fit_transform(df_verse_complete_dict_nrmzl['verse'])
tfidf_matrix_nahj_char = vectorizer_nahj_char.fit_transform(df_nahj_complete_dict_nrmzl['verse'])
tfidf_matrix_sahife_char = vectorizer_sahife_char.fit_transform(df_sahife_complete_dict_nrmzl['verse'])
#

# Print shapes of tfidf_matrices
#
print("Shape of tfidf_matrix(Word-Level) of quran is: ",tfidf_matrix_quran.shape)
print("Shape of tfidf_matrix(Word-Level) of Nahj is: ",tfidf_matrix_nahj.shape)
print("Shape of tfidf_matrix(Word-Level) of Sahife is: ",tfidf_matrix_sahife.shape)
print("Shape of tfidf_matrix of quran is: ",tfidf_matrix_quran_char.shape)
print("Shape of tfidf_matrix of Nahj is: ",tfidf_matrix_nahj_char.shape)
print("Shape of tfidf_matrix of Sahife is: ",tfidf_matrix_sahife_char.shape)
#

Shape of tfidf_matrix(Word-Level) of quran is: (6236, 14658)
Shape of tfidf_matrix(Word-Level) of Nahj is: (800, 21094)
Shape of tfidf_matrix(Word-Level) of Sahife is: (924, 7711)
Shape of tfidf_matrix of quran is: (6236, 40)
Shape of tfidf_matrix of Nahj is: (800, 36)
Shape of tfidf_matrix of Sahife is: (924, 47)
```

سپس یک لیست از زوج (ایندکس:آیه) از هر سه کتاب استخراج می‌کنیم زیرا در ادامه نیاز به جستجو آیه بر اساس ایندکس یا برعکس می‌شود. در نهایت در بخش بعد ورودی مدل را مشخص کرده تا در ادامه از هر کتاب، ۱۰ آیه/جمله مشابه از لحاظ معیار فاصله فضای تعبیه ساخته شده توسط وکتورایزر TF-IDF نشان داده می‌شود. برای محاسبه MRR به جهت مقایسه این روش شباهت‌یابی و دو روش بعدی، ترتیب ۱۰ آیه/جمله مشابه را بهم ریخته و به دو نفر اعضای تیم که این خروجی‌ها را ندیده‌اند، داده می‌شود تا به صورت دستی مشابه‌ترین آیه/جمله مد نظرشان را اعلام کنند.

## یافتن عبارات مشابه در متون دینی با استفاده از جاگذاری کلمات

تکلیف سری سوم

ورودی که برای روش اول یعنی TF-IDF آزمایش شد، یکی از آیه‌های قرآن یعنی آیه ۴ سوره ناس: «مِنْ شَرِّ الْوَسْوَاسِ الْخَنَّاسِ» بود. توجه شود که چون یکی از آیه‌های قرآن را به عنوان ورودی مشخص کردیم، طبیعی است که برترین و مشابه‌ترین آیه قرآنی با این آیه، خودش باشد لذا برای قرآن، ما ۹ آیه مشابه بعدی را برای محاسبه دستی استفاده کردیم. ۹ آیه برتر در شباهت به ترتیب عبارتند از:

۱- مِنْ شَرِّ مَا خَلَقَ

۲- وَمِنْ شَرِّ النَّفَّاثَاتِ فِي الْعُقَدِ

۳- وَمِنْ شَرِّ غَاسِقٍ إِذَا وَقَبَ

۴- وَمِنْ شَرِّ حَاسِدٍ إِذَا حَسَدَ

۵- إِنَّ شَرَّ الدَّوَابِّ عِنْدَ اللَّهِ الَّذِينَ كَفَرُوا فَهُمْ لَا يُؤْمِنُونَ

۶- إِنَّ شَرَّ الدَّوَابِّ عِنْدَ اللَّهِ الصُّمُّ الْبُكْمُ الَّذِينَ لَا يَعْقِلُونَ

۷- إِنَّ الَّذِينَ كَفَرُوا مِنْ أَهْلِ الْكِتَابِ وَالْمُشْرِكِينَ فِي نَارِ جَهَنَّمَ خَالِدِينَ فِيهَا ۖ أُولَٰئِكَ هُمْ شَرُّ الْبَرِيَّةِ

۸- الَّذِينَ يُحْشَرُونَ عَلَىٰ وُجُوهِهِمْ إِلَىٰ جَهَنَّمَ أُولَٰئِكَ شَرُّ مَكَانًا وَأَضَلُّ سَبِيلًا

۹- فَوَقَاهُمُ اللَّهُ شَرَّ ذَٰلِكَ الْيَوْمِ وَلَقَّاهُمْ نَضْرَةً وَسُرُورًا

نفر اول آیه شماره ۲ و نفر دوم آیه شماره ۴ را به عنوان مشابه‌ترین آیه برگزیدند. بنابراین داریم:

$$MRR_{QuranTF-IDF} = \frac{1}{2} \sum_{i=1}^2 \frac{1}{Rank_{choice}^i} = \frac{1}{2} \left[ \frac{1}{2} + \frac{1}{4} \right] = 0.375$$

به همین ترتیب آیه/جمله‌های برتر در نهج البلاغه عبارت بودند از:

۱- وَإِنَّ عَمَلَكَ لَيْسَ لَكَ بِطُعْمَةٍ وَلَكِنَّهُ فِي عُنُقِكَ أَمَانَةٌ وَأَنْتَ مُسْتَرْعَى لِمَنْ فَوْقَكَ لَيْسَ لَكَ أَنْ تَفْتَتَ فِي رَعِيَّةٍ وَلَا تُخَاطِرَ إِلَّا بِوَيْفَقَةٍ وَفِي يَدَيْكَ مَالٌ مِنْ مَالِ اللَّهِ عَزَّ وَجَلَّ وَأَنْتَ مِنْ خَزَائِنِهِ حَتَّى تُسَلِّمَهُ إِلَيَّ وَلَعَلِّي أَلَا أَكُونُ شَرًّا وَلَاتِكَ لَكَ وَالسَّلَامُ

۲- فَأَخْفِضْ لَهُمْ جَنَاحَكَ وَارْلِنْ لَهُمْ جَانِبَكَ وَابْسُطْ لَهُمْ وَجْهَكَ وَاسِ بَيْنَهُمْ فِي اللَّحْظَةِ وَالنَّظَرَةِ حَتَّى لَا يَطْمَعَ الْعُظَمَاءُ فِي حَيْفِكَ لَهُمْ وَلَا يَبْأَسَ الضَّعَفَاءُ مِنْ عَدْلِكَ عَلَيْهِمْ فَإِنَّ اللَّهَ تَعَالَى يُسَائِلُكُمْ مَعَشَرَ عِبَادِهِ عَنِ الصَّغِيرَةِ مِنْ أَعْمَالِكُمْ وَالْكَبِيرَةِ وَالظَّاهِرَةِ وَالْمُسْتَوْرَةِ فَإِنْ يُعَذِّبْ فَأَنْتُمْ أَظْلَمُ وَإِنْ يَعْفُ فَهُوَ أَكْرَمُ وَاعْلَمُوا عِبَادَ اللَّهِ أَنَّ الْمُتَّقِينَ ذَهَبُوا بِعَاجِلِ الدُّنْيَا وَآجِلِ



الْآخِرَةِ فَشَارَكُوا أَهْلَ الدُّنْيَا فِي دُنْيَاهُمْ وَلَمْ يُشَارِكُوا أَهْلَ الدُّنْيَا فِي آخِرَتِهِمْ سَكَنُوا الدُّنْيَا بِأَفْضَلِ مَا سَكَنَتْ وَ أَكَلُوا بِأَفْضَلِ مَا أَكَلَتْ فَحَظُّوا مِنَ الدُّنْيَا بِمَا حَظَّى بِهِ الْمُتَرَفُّونَ وَ أَخَذُوا مِنْهَا مَا أَخَذَهُ الْجَبَّارَةُ الْمُتَكَبِّرُونَ ثُمَّ انْقَلَبُوا عَنْهَا بِالزَّادِ الْمُبْلَغِ وَ الْمَتَجَرِّ الرَّابِحِ أَصَابُوا لَذَّةَ زُهْدِ الدُّنْيَا فِي دُنْيَاهُمْ وَ تَيَقَّنُوا أَنَّهُمْ حَيْرَانُ اللَّهِ غَدًا فِي آخِرَتِهِمْ لَا تُرَدُّ لَهُمْ دَعْوَةٌ وَ لَا يَنْقُصُ لَهُمْ نَصِيبٌ مِنْ لَذَّةٍ فَاحْذَرُوا عِبَادَ اللَّهِ الْمَوْتَ وَ قُرْبَهُ وَ أَعْدُوا لَهُ عَدُوَّهُ فَإِنَّهُ يَأْتِي بِأَمْرِ عَظِيمٍ وَ خَطْبٍ جَلِيلٍ بِخَيْرٍ لَا يَكُونُ مَعَهُ شَرٌّ أَبَدًا أَوْ شَرٌّ لَا يَكُونُ مَعَهُ خَيْرٌ أَبَدًا فَمَنْ أَقْرَبُ إِلَى الْجَنَّةِ مِنْ عَامِلِهَا وَ مَنْ أَقْرَبُ إِلَى النَّارِ مِنْ عَامِلِهَا وَ أَنْتُمْ طُرْدَاءُ الْمَوْتِ إِنْ أَقَمْتُمْ لَهُ أَخَذَكُمْ وَ إِنْ فَرَرْتُمْ مِنْهُ أَدْرَكَكُمْ وَ هُوَ أَلْزَمُ لَكُمْ مِنْ ظِلِّكُمْ الْمَوْتُ مَعْقُودٌ بِنَوَاصِيكُمْ وَ الدُّنْيَا تَطْوَى مِنْ خَلْفِكُمْ فَاحْذَرُوا نَاراً قَعْرُهَا بَعِيدٌ وَ حَرُّهَا شَدِيدٌ وَ عَذَابُهَا جَدِيدٌ دَارٌ لَيْسَ فِيهَا رَحْمَةٌ وَ لَا تَسْمَعُ فِيهَا دَعْوَةٌ وَ لَا تُفَرِّجُ فِيهَا كُرْبَةٌ وَ إِنْ اسْتَطَعْتُمْ أَنْ يَشْتَدَّ خَوْفُكُمْ مِنَ اللَّهِ وَ أَنْ يَحْسُنَ ظَنُّكُمْ بِهِ فَاجْمَعُوا بَيْنَهُمَا فَإِنَّ الْعَبْدَ إِنَّمَا يَكُونُ حَسَنُ ظَنِّهِ بِرَبِّهِ عَلَى قَدْرِ خَوْفِهِ مِنْ رَبِّهِ وَ إِنْ أَحْسَنَ النَّاسُ ظَنًّا بِاللَّهِ أَشَدَّهُمْ خَوْفًا لِلَّهِ وَ أَعْلَمَ يَا مُحَمَّدُ بْنُ أَبِي بَكْرٍ أَنِّي قَدْ وَلَّيْتُكَ أَعْظَمَ أَجْنَادِي فِي نَفْسِي أَهْلَ مِصْرَ فَأَنْتَ مُحَقَّقٌ أَنْ تُخَالَفَ عَلَى نَفْسِكَ وَ أَنْ تُنَافِحَ عَنْ دِينِكَ وَ لَوْ لَمْ يَكُنْ لَكَ إِلَّا سَاعَةٌ مِنَ الدَّهْرِ وَ لَا تُسَخِّطِ اللَّهَ بِرِضَا أَحَدٍ مِنْ خَلْقِهِ فَإِنَّ فِي اللَّهِ خَلْفًا مِنْ غَيْرِهِ وَ لَيْسَ مِنَ اللَّهِ خَلْفٌ فِي غَيْرِهِ صَلِّ الصَّلَاةَ لَوَقْتِهَا الْمُؤَقَّتِ لَهَا وَ لَا تُعْجَلْ وَقْتُهَا لِغَرَاغٍ وَ لَا تُؤَخِّرْهَا عَنْ وَقْتِهَا لِاسْتِغَالٍ وَ أَعْلَمْ أَنَّ كُلَّ شَيْءٍ مِنْ عَمَلِكَ تَبَعَ لِصَلَاتِكَ وَ مِنْهُفَانُهُ لَا سَوَاءَ إِمَامُ الْهُدَى وَ إِمَامُ الرَّدَى وَ وَلِيَّ النَّبِيِّ وَ عَدُوُّ النَّبِيِّ وَ لَقَدْ قَالَ لِي رَسُولُ اللَّهِ -ص- إِنِّي لَا أَخَافُ عَلَى أُمَّتِي مُؤْمِنًا وَ لَا مُشْرِكًا أَمَّا الْمُؤْمِنُ فَيَمْنَعُهُ اللَّهُ بِإِيمَانِهِ وَ أَمَّا الْمُشْرِكُ فَيَقْمَعُهُ اللَّهُ بِشِرْكِهِ وَ لَكِنِّي أَخَافُ عَلَيْكُمْ كُلَّ مُنَافِقٍ الْجَنَانِ عَالِمِ اللِّسَانِ يَقُولُ مَا تَعْرِفُونَ وَ يَفْعَلُ مَا تُنْكِرُونَ

۳- وَ قَالَ عَ شَرُّ الْإِخْوَانِ مَنْ تُكَلِّفَ لَهُ

۴- وَ قَالَ عَ يَأْتِي عَلَى النَّاسِ زَمَانٌ لَا يَبْقَى فِيهِمْ مِنَ الْقُرْآنِ إِلَّا رَسْمُهُ وَ مِنَ الْإِسْلَامِ إِلَّا اسْمُهُ وَ مَسَاجِدُهُمْ يَوْمِئِذٍ عَامِرَةٌ مِنَ الْبِنَاءِ خَرَابٌ مِنَ الْهُدَى سُكَّانُهَا وَ عِمَارُهَا شَرُّ أَهْلِ الْأَرْضِ مِنْهُمْ تَخْرُجُ الْفِتْنَةُ وَ إِلَيْهِمْ تَأْوِي الْخَطِيئَةُ يَرُدُّونَ مَنْ شَذَّ عَنْهَا فِيهَا وَ يَسُوقُونَ مَنْ تَأَخَّرَ عَنْهَا إِلَيْهَا يَقُولُ اللَّهُ سُبْحَانَهُ فَبِي حَلَفْتُ لَا أَبْعَثَنَّ عَلَى أَوْلَيْكَ فِتْنَةً تَتْرُكُ الْحَلِيمَ فِيهَا حَيْرَانَ وَ قَدْ فَعَلَ وَ نَحْنُ نَسْتَقِيلُ اللَّهَ عَثْرَةَ الْعَفْلَةِ

۵- وَ قَالَ عَ فَاعِلُ الْخَيْرِ خَيْرٌ مِنْهُ وَ فَاعِلُ الشَّرِّ شَرٌّ مِنْهُ

۶- وَ قَالَ عَ الْمَرْأَةُ شَرُّ كُلِّهَا وَ شَرٌّ مَا فِيهَا أَنَّهُ لَا بُدَّ مِنْهَا

۷- وَ قَالَ عَ مَا خَيْرٌ بِخَيْرٍ بَعْدَهُ النَّارُ وَ مَا شَرٌّ بِشَرٍّ بَعْدَهُ الْجَنَّةُ وَ كُلُّ نَعِيمٍ دُونَ الْجَنَّةِ فَهُوَ مُحَقَّقٌ وَ كُلُّ بَلَاءٍ دُونَ النَّارِ عَافِيَةٌ

۸- إِنْ اللَّهَ بَعَثَ مُحَمَّدًا -ص- نَذِيرًا لِلْعَالَمِينَ وَ أَمِينًا عَلَى التَّنْزِيلِ وَ أَنْتُمْ مَعَشَرَ الْعَرَبِ عَلَى شَرِّ دِينٍ وَ فِي شَرِّ دَارٍ مُنِيخُونَ بَيْنَ حِجَارَةٍ خُشْنٍ وَ حَيَاتٍ صُمٍّ تَشْرَبُونَ الْكَدِرَ وَ تَأْكُلُونَ الْجَشِبَ وَ تَسْفِكُونَ دِمَاءَكُمْ وَ تَقْطَعُونَ أَرْحَامَكُمْ الْأَصْنَامُ فِيكُمْ مَنْصُوبَةٌ وَ الْإِثَامُ بِكُمْ مَعْصُوبَةٌ فَتَنْظَرْتُ فَإِذَا لَيْسَ لِي مُعِينٌ إِلَّا أَهْلُ بَيْتِي فَضَنَنْتُ بِهِمْ عَنِ الْمَوْتِ وَ أَغْضَيْتُ عَلَى الْقَذَى وَ شَرِبْتُ عَلَى

الشَّجَا وَ صَبْرَتْ عَلَى أَخَذِ الْكَظْمِ وَ عَلَى أَمْرٍ مِنْ طَعْمِ الْعَلَقَمِ وَمِنْهَا وَ لَمْ يُبَايِعْ حَتَّى شَرَطَ أَنْ يُؤْتِيَهُ عَلَى الْبَيْعَةِ ثَمَنًا فَلَا ظَفِرَتْ يَدُ الْبَائِعِ وَ خَزَيْتَ أَمَانَةَ الْمُبْتَاعِ فَخَذُوا لِلْحَرْبِ أَهْبَتَهَا وَ أَعَدُّوا لَهَا عُدَّتَهَا فَقَدْ شَبَّ لَهَا وَ عَلَا سَنَاهَا وَ اسْتَشَعَرُوا الصَّبْرَ فَإِنَّهُ أَدْعَى إِلَى النَّصْرِ

۹- فَإِنَّكَ قَدْ جَعَلْتَ دِينَكَ تَبَعًا لِدُنْيَا امْرِئٍ ظَاهِرٍ غَيْهٍ مَهْتُوكٍ سِتْرَهُ يَشِينُ الْكَرِيمَ بِمَجْلِسِهِ وَ يُسَفِّهُ الْحَلِيمَ بِخِلَاطَتِهِ فَاتَّبَعْتَ أَثَرَهُ وَ طَلَبْتَ فَضْلَهُ اتَّبَعَ الْكَلْبُ لِلضَّرْعَامِ يَلُودُ بِمَخَالِبِهِ وَ يَنْتَظِرُ مَا يُلْقَى إِلَيْهِ مِنْ فَضْلٍ فَرِيَسْتِهِ فَأَذْهَبْتَ دُنْيَاكَ وَ آخَرْتَكَ وَ لَوْ بِالْحَقِّ أَخَذْتَ أَدْرَكَتَ مَا طَلَبْتَ فَإِنْ يُمْكِنِي اللَّهُ مِنْكَ وَ مِنْ ابْنِ أَبِي سُفْيَانَ أَجْزِكُمَا بِمَا قَدَّمْتُمَا وَ إِنْ تُعْجِزَا وَ تَبْقَيَا فَمَا أَمَامَكُمَا شَرٌّ لَكُمَا وَ السَّلَامُ

۱۰- صَابَكُمْ حَاصِبٌ وَ لَا بَقِيَ مِنْكُمْ أَثَرٌ أَوْ بَعْدَ إِيمَانِي بِاللَّهِ وَ جِهَادِي مَعَ رَسُولِ اللَّهِ-ص أَشْهَدُ عَلَى نَفْسِي بِالْكَفْرِ لَقَدْ ضَلَلْتُ إِذَا وَ مَا أَنَا مِنَ الْمُهْتَدِينَ فَأَوْبُوا شَرَّ مَا بٍ وَ ارْجِعُوا عَلَى أَثَرِ الْأَعْقَابِ أَمَا إِنَّكُمْ سَتَلْقَوْنَ بَعْدِي ذُلًّا شَامِلًا وَ سَيْفًا قَاطِعًا وَ أَثَرَهُ يَتَخَذُهَا الظَّالِمُونَ فِيكُمْ سُنَّةً

نفر اول آیه شماره ۸ و نفر دوم آیه شماره ۶ را به عنوان مشابه ترین آیه برگزیدند. بنابراین داریم:

$$MRR_{Nahj_{TF-IDF}} = \frac{1}{2} \sum_{i=1}^2 \frac{1}{Rank_{choice}^i} = \frac{1}{2} \left[ \frac{1}{8} + \frac{1}{6} \right] = 0.146$$

به همین ترتیب آیه/جمله های برتر در صحیفه سجاده عبارت بودند از:

۱- فَسُبْحَانَكَ! مَا أَبِينَ كَرَمَكَ فِي مُعَامَلَةٍ مَنْ أَطَاعَكَ أَوْ عَصَاكَ: تَشْكُرُ لِلْمُطِيعِ مَا أَنْتَ تَوَكَّلْتَهُ لَهُ، وَ تُدْمِلِي لِلْعَاصِي فِيمَا تَمْلِكُ مُعَاجَلَتَهُ فِيهِ.

۲- وَ أَجْزَلُ لَنَا فِيهِ مِنَ الْحَسَنَاتِ، وَ أَخْلَنَّا فِيهِ مِنَ السَّيِّئَاتِ، وَ أَمَلْنَا لَنَا مَا بَيْنَ طَرَفَيْهِ حَمْدًا وَ شُكْرًا وَ أَجْرًا وَ ذُخْرًا وَ فَضْلًا وَ إِحْسَانًا.

۳- وَ أَمْنُنْ عَلَيَّ بِالصَّحَّةِ وَ الْأَمْنِ وَ السَّلَامَةِ فِي دِينِي وَ بَدَنِي، وَ الْبَصِيرَةِ فِي قَلْبِي، وَ النَّفَازِ فِي أُمُورِي، وَ الْخَشْيَةِ لَكَ، وَ الْخَوْفِ مِنْكَ، وَ الْقُوَّةِ عَلَى مَا أَمَرْتَنِي بِهِ مِنْ طَاعَتِكَ، وَ الْاجْتِنَابِ لِمَا نَهَيْتَنِي إِذَا أَنْصَرَفَ مِنْ صَلَاتِهِ قَامَ قَائِمًا ثُمَّ اسْتَقْبَلَ الْقِبْلَةَ، وَ فِي يَوْمِ الْجُمُعَةِ، فَقَالَ: عَنْهُ مِنْ مَعْصِيَتِكَ.

۴- وَ أَنْ نَتَقَرَّبَ إِلَيْكَ فِيهِ مِنَ الْأَعْمَالِ الزَّكَايَةِ بِمَا تُطَهِّرُنَا بِهِ مِنَ الذُّنُوبِ، وَ تَعْصِمُنَا فِيهِ مِمَّا نَسْتَأْنِفُ مِنَ الْعُيُوبِ، حَتَّى لَا يُورِدَ عَلَيْكَ أَحَدٌ مِنْ مَلَائِكَتِكَ إِلَّا دُونَ مَا نُورِدُ مِنْ أَبْوَابِ الطَّاعَةِ لَكَ، وَ أَنْوَاعِ الْقُرْبَةِ إِلَيْكَ.

۵- وَ اسْتَخْلِكَ مِنْ ذُنُوبِي مَا قَدْ بَهَظَنِي حَمْلُهُ، وَ اسْتَعِينُ بِكَ عَلَى مَا قَدْ فَدَحَنِي ثِقَلُهُ.

۶- وَ لَا تَذَرْنِي فِي طُغْيَانِي عَامِيًا، وَ لَا فِي غَمْرَتِي سَاهِيًا حَتَّى حِينٍ، وَ لَا تَجْعَلَنِي عِظَةً لِمَنْ اتَّعَظَ، وَ لَا نَكَالًا لِمَنْ اعْتَبَرَ، وَ لَا فِتْنَةً لِمَنْ نَظَرَ، وَ لَا تَمْكُرْ بِي فِيمَنْ تَمْكُرُ بِهِ، وَ لَا تَسْتَبْدِلْ بِي غَيْرِي، وَ لَا تُغَيِّرْ لِي اسْمًا، وَ لَا تُبَدِّلْ لِي جِسْمًا، وَ لَا تَتَّخِذْنِي هُزُوءًا لِخَلْقِكَ، وَ لَا سُخْرِيًّا لَكَ، وَ لَا تَبْعًا إِلَّا لِمَرْضَاتِكَ، وَ لَا مُمْتَهَنًا إِلَّا بِالْإِنْتِقَامِ لَكَ.

۷- اللَّهُمَّ اجْعَلْ مَا يُلْقَى الشَّيْطَانُ فِي رُوعِي مِنَ التَّمَنَّى وَ التَّظَنِّي وَ الْحَسَدِ ذِكْرًا لِعَظَمَتِكَ، وَ تَفَكُّرًا فِي قُدْرَتِكَ، وَ تَذِيبًا عَلَى عَدُوِّكَ، وَ مَا أَجْرَى عَلَى لِسَانِي مِنْ لَفْظَةٍ فُحْشٍ أَوْ هُجْرٍ أَوْ شَتْمٍ عَرَضٍ أَوْ شَهَادَةٍ بَاطِلٍ أَوْ اغْتِيَابٍ مُؤْمِنٍ غَائِبٍ أَوْ سَبِّ حَاضِرٍ وَ مَا أَشَبَهُ ذَلِكَ نَطْقًا بِالْحَمْدِ لَكَ، وَ إِغْرَاقًا فِي الثَّنَاءِ عَلَيْكَ، وَ ذَهَابًا فِي تَمْجِيدِكَ، وَ شُكْرًا لِنِعْمَتِكَ، وَ اعْتِرَافًا بِإِحْسَانِكَ، وَ إِحْصَاءً لِمِنَّكَ.

۸- وَ نَعُوذُ بِكَ أَنْ نَنْطَوِيَ عَلَى غِشٍّ أَحَدٍ، وَ أَنْ نُعْجِبَ بِأَعْمَالِنَا، وَ نَمُدَّ فِي آمَالِنَا

۹- اللَّهُمَّ صَلِّ عَلَى مُحَمَّدٍ وَ آلِهِ، وَ اكْفِنِي مَا يَشْغَلُنِي الْإِهْتِمَامُ بِهِ، وَ اسْتَعْمِلْنِي بِمَا تَسْأَلُنِي غَدًا عَنْهُ، وَ اسْتَفْرِغْ أَيْامِي فِيمَا خَلَقْتَنِي لَهُ، وَ اغْنِنِي وَ أَوْسِعْ عَلَيَّ فِي رِزْقِكَ، وَ لَا تَفْتِنَنِي بِالنَّظَرِ، وَ اعِزَّنِي وَ لَا تَبْلِيَنِي بِالْكِبَرِ، وَ عَبْدُنِي لَكَ وَ لَا تُفْسِدْ عِبَادَتِي بِالْعُجْبِ، وَ أَجِرْ لِلنَّاسِ عَلَى يَدِي الْخَيْرَ وَ لَا تَمَحِّقْهُ بِالْمَنِّ، وَ هَبْ لِي مَعَالِيَ الْأَخْلَاقِ، وَ اعْصِمْنِي مِنَ الْفَخْرِ.

۱۰- حَتَّى اسْتَتَبَ لَهُ مَا حَاوَلَ فِي أَعْدَائِكَ

نفر اول آیه شماره ۷ و نفر دوم آیه شماره ۲ را به عنوان مشابه ترین آیه برگزیدند. بنابراین داریم:

$$\text{MRR}_{\text{SahifaTF-IDF}} = \frac{1}{2} \sum_{i=1}^2 \frac{1}{\text{Rank}_{\text{choice}}^i} = \frac{1}{2} \left[ \frac{1}{7} + \frac{1}{2} \right] = 0.321$$

بنابراین اگر روش TF-IDF Word-Level را بخواهیم با یک MRR نشان دهیم، با میانگین گیری از این سه عدد بدست آمده داریم:

## یافتن عبارات مشابه در متون دینی با استفاده از جاگذاری کلمات

تکلیف سری سوم

$$MRR_{TF-IDF} = \frac{0.375 + 0.146 + 0.321}{3} = 28\%$$

به همین ترتیب برای حالتی که روی فضای کاراکترها تعبیه‌سازی کنیم، نفر اول برای هر ۳ کتب، آیه‌های ۳ و ۵ و ۹ و نفر دوم آیه‌های ۴ و ۳ و ۳ را انتخاب کردند لذا داریم:

$$MRR_{Quran_{TF-IDF}} = \frac{1}{2} \sum_{i=1}^2 \frac{1}{Rank_{choice}^i} = \frac{1}{2} \left[ \frac{1}{3} + \frac{1}{4} \right] = 0.291$$

$$MRR_{Nahj_{TF-IDF}} = \frac{1}{2} \sum_{i=1}^2 \frac{1}{Rank_{choice}^i} = \frac{1}{2} \left[ \frac{1}{5} + \frac{1}{3} \right] = 0.266$$

$$MRR_{Sahifa_{TF-IDF}} = \frac{1}{2} \sum_{i=1}^2 \frac{1}{Rank_{choice}^i} = \frac{1}{2} \left[ \frac{1}{9} + \frac{1}{3} \right] = 0.222$$

پس داریم:

$$MRR_{TF-IDF} = \frac{0.291 + 0.266 + 0.222}{3} = 26\%$$

## روش دوم – استفاده از FASTTEXT

در روش دوم برای محاسبه شباهت بین جملات از مدل FastText استفاده کردیم. این مدل توسط شرکت فیسبوک توسعه داده شده و هم اکنون از پرستفاده‌ترین ترنسفورمرهای دنیای پردازش زبان طبیعی است. در این روش نیز از همان موارد پیش‌پردازش ذکر شده در بخش اول استفاده شده است. برای بهره‌گیری از قابلیت‌های این ابزار از کتابخانه Gensim استفاده کردیم که یک کتابخانه برای Topic Modeling می‌باشد. همچنین برای تولید وکتورها و تعبیه آنها از مدل پیش‌آموزش داده شده عربی فیسبوک استفاده کردیم.

Load pretrained arabic vector (downloaded from fasttext data storage)

```
[ ] 1 model = KeyedVectors.load_word2vec_format(datapath(f'{vectorModelDir}cc.ar.300.vec'), binary=False)
```

بعد از لود کردن مدل نوبت به تولید وکتور برای کلمات موجود مجموعه داده‌ها می‌رسد (قرآن کریم، نهج البلاغه و صحیفه سجادیه). برای این کار توکن‌های هر آیه/جمله را به مدل داده تا برای آنها وکتور تولید کند. این کار با استفاده از تابع زیر صورت می‌گیرد:

```
1
2 """
3 get vector for each token from pretrained arabic model
4 if each token doesnt exist in model we assign 0 to its vector
5 @param tokens an array of sentence word
6 @return an array of tuples with format (token, vector)
7 """
8 def w2v4corpus(token, models):
9     w2vToken = token
10    k=0
11    for i in range(len(token)):
12        for j in range(len(token[i])):
13            word = token[i][j]
14            models = model
15            if word in models:
16                w2vToken[i][j] = (word, models[word])
17            else:
18                # list of words that doesnt exist in the vector
19                print(word)
20                k=k+1
21                w2vToken[i][j] = (word, 0)
22    print(f'number of words that doesnt exist in the pretrained fasttext model: {k}')
23    return w2vToken
```

بعضی از لغات در مجموعه داده ما بودند که در مدل فیسبوک برای آن تعبیه‌سازی‌ای وجود نداشت. البته بسیاری از آنها به دلیل نوشتار نادرست با این مشکل روبه‌رو شده بودند (مانند چسبیدن حرف «و» به کلمه بعدی) که باعث شد ما آنها را

# یافتن عبارات مشابه در متون دینی با استفاده از جاگذاری کلمات

تکلیف سری سوم

نادیده بگیریم. متاسفانه مدل FastText مانند مدل BERT از توکنایزر پیشرفته و مختص به خود بهره نمی‌برد و استفاده از ابزارهای متفاوت توکنایزیشن نیز برای ما دردی را دوا نکرد به همین دلیل تصمیم به چشم پوشی از این کلمات گرفتیم. بعد از تشکیل وکتور برای کلمات حاضر در آیات/جملات نوبت به تشکیل وکتور برای آیات/جملات می‌رسد. برای این کار می‌توانستیم از دو روش استفاده کنیم. اول از طریق API های خود FastText و دوم ساخت آن به صورت دستی. برای استفاده از FastText نیاز داشتیم تا یک مدل Doc2Vec برای زبان عربی در دسترس داشته باشیم که بعد از جست و جوی فراوان پیدا نشد (تقریباً فقط برای زبان انگلیسی موجود بود). پس به سراغ راه دوم، یعنی ساخت وکتور به صورت دستی، رفتیم. برای این کار وکتور جمله را میانگین وکتور کلمات سازنده آن در نظر گرفتیم. در قطعه کد زیر توابع مورد نیاز برای این مورد آورده شده است:

```
1 # calculate avg of word vectors in sentence
2 def calculateSentenceVector(wordVector):
3     return np.mean( np.array(wordVector), axis=0 )
4
5 # convert doc to vector
6 def doc2vec(corpus):
7     temp = []
8     for i in range(len(corpus)):
9         temp.append(calculateSentenceVector([l[1] for l in corpus[i]]))
10    return temp
```

پس از این مرحله ما وکتور را برای تمام آیات و جملات موجود در مجموعه داده در دسترس داریم. پس نوبت به کوئری زدن و یافتن موارد مشابه است. برای کوئری ورودی نیز تمام مراحل پیش‌پردازش و وکتور کردن که در این بخش ذکر کردیم صادق است. بعد از تشکیل وکتور برای جمله ورودی، با استفاده از محاسبه کسینوس زاویه بین وکتور جمله ورودی (کوئری) و تک تک جملات موجود در مجموعه داده، موارد محاسبه شده به صورت نزولی مرتب شده و ۱۰ تایی اول به عنوان ۱۰ جمله با بالاترین میزان تشابه انتخاب می‌شوند.

```
1 """
2 calculate similarity between input query and all words in vecList using cosine similarity
3 @param query the input query vector
4 @param vecList a list of vectors(words)
5 @return a list of similarity rank
6 """
7
8 def calculateSimilarity(qurey, vecList):
9     temp=[]
10    k=0
11    for i in vecList[1]:
12        print(k)
13        k = k+1
14        temp.append(cosine_similarity([qurey], [i]))
15    return temp
16
```

# یافتن عبارات مشابه در متون دینی با استفاده از جاگذاری کلمات

تکلیف سری سوم

تابع زیر وظیفه یافتن مشابه‌ترین‌ها را بر عهده دارد. این امکان وجود دارد که در ورودی تابع مقرر کنیم تا در کدام متون دنبال جملات مشابه بگردیم.

```
"""
return similarity matrix
@param quran if true search in quranic verse to find similar verse
@param nahj if true search in nahj ol balaghe verse to find similar verse
@param sahifa if true search in sahifa verse to find similar verse
"""
def mostSimilar(queryVector, quran=True, nahj=False, sahifa=False):
    similarVector=[]
    if quran:
        similarVector.append(calculateSimilarity(queryVector, verse_complete_dict_nrmlze))
    if nahj:
        similarVector.append(calculateSimilarity(queryVector, nahj_complete_dict_nrmlze))
    if sahifa:
        similarVector.append(calculateSimilarity(queryVector, sahifa_complete_dict_nrmlze))
    return similarVector
```

```
1 # a list for similarity between input query and all sentence
2 similarity = mostSimilar(query, quran=True)

[ ] 1 #sort similarity and select top 10
     2 sortedSimilarity = sorted(similarity, key=lambda x: x[1])
     3 sortedSimilarity[:10]
```

برای ارزیابی مدل طبق روش MRR ما یک آیه را به عنوان ورودی به مدل دادیم و سپس خروجی زیر را به عنوان ۱۰ آیه با بیشترین تشابه دریافت کردیم (بدیهی است که شبیه‌ترین عبارت به یک آیه، خود آیه است):

ورودی:

[ 'وَإِذْ قُلْنَا لِلْمَلَائِكَةِ اسْجُدُوا لِآدَمَ فَسَجَدُوا إِلَّا إِبْلِيسَ أَبَىٰ وَاسْتَكْبَرَ وَكَانَ مِنَ الْكَافِرِينَ' ]

خروجی:

# یافتن عبارات مشابه در متون دینی با استفاده از جاگذاری کلمات

تکلیف سری سوم

وَإِذْ قُلْنَا لِلْمَلَائِكَةِ اسْجُدُوا لِآدَمَ فَسَجَدُوا إِلَّا إِبْلِيسَ أَبَىٰ وَاسْتَكْبَرَ وَكَانَ مِنَ الْكَافِرِينَ
وَلَقَدْ خَلَقْنَاكُمْ ثُمَّ صَوَّرْنَاكُمْ ثُمَّ قُلْنَا لِلْمَلَائِكَةِ اسْجُدُوا لِآدَمَ فَسَجَدُوا إِلَّا إِبْلِيسَ لَمْ يَكُنْ مِنَ السَّاجِدِينَ
وَإِذْ قُلْنَا لِلْمَلَائِكَةِ اسْجُدُوا لِآدَمَ فَسَجَدُوا إِلَّا إِبْلِيسَ قَالَ أَأَسْجُدُ لِمَنْ خَلَقْتُ طِينًا
وَإِذْ قُلْنَا لِلْمَلَائِكَةِ اسْجُدُوا لِآدَمَ فَسَجَدُوا إِلَّا إِبْلِيسَ أَبَىٰ
﴿وَإِذَا قِيلَ لَهُمْ اسْجُدُوا لِلرَّحْمَنِ قَالُوا وَمَا الرَّحْمَنُ أَنَسْجُدُ لِمَا تَأْمُرُنَا وَزَادَهُمْ نُفُورًا
أَمْ يَقُولُونَ إِنَّا بِرَأْيِهِمْ نَحْنُ الْغَاثِلُونَ إِنَّا أَنزَلْنَاهُ فَاذْكُرُونَا أَنَّهُ لَا إِلَهَ إِلَّا اللَّهُ وَبِهِ نَعِيبُونَ أَلَمْ يَكُنْ لَهُ الْإِلَهَ الْأُولَىٰ فَلَمْ يُهَيِّئْ لَهُ شَرَفًا فَوُضِعَ فِي الْآثَانِ إِنَّا أَنزَلْنَاهُ فِي الْقُرْآنِ وَإِن تَوَلَّوْا فَإِنَّمَا تَوَلَّوْا الْكُفْرَ
إِلَّا إِبْلِيسَ اسْتَكْبَرَ وَكَانَ مِنَ الْكَافِرِينَ
إِذْ قَالَ يُوسُفُ لِأَبِيهِ يَا أَبَتِ إِنِّي رَأَيْتُ أَحَدَ عَشَرَ كَوْكَبًا وَالشَّمْسَ وَالْقَمَرَ رَأَيْتُهُمْ لِي سَاجِدِينَ
فَلَمَّا رَأَىٰ الشَّمْسُ بَازِعَهُ قَالَ هَذَا رَبِّي هَذَا أَكْبَرُ فَلَمَّا أَفَلَتْ قَالَ يَا قَوْمِ إِنِّي بَرِيءٌ مِّمَّا تُشْرِكُونَ
قَالَ الْمَلَأُ الَّذِينَ اسْتَكْبَرُوا مِنْ قَوْمِهِ لِلَّذِينَ اسْتُضْعِفُوا لِمَنْ آمَنَ مِنْهُمْ أَتَعْلَمُونَ إِنَّ صَالِحًا مَّرْسَلًا مِنْ رَبِّهِ قَالُوا إِنَّا بِمَا أُرْسِلَ بِهِ مُؤْمِنُونَ

سپس برای ارزیابی مدل (نابود شدیم پای این قسمت تا جملات راحت پیدا کنیم که راحت تر بشه مقایسه کرد!!!)، به مانند بخش TF-IDF این ۱۰ آیه را به صورت نامرتب به ۲ عضو دیگر دادیم تا مشابه ترین را مشخص کنند که هر دو نفر جایگاه ۱ را حدس زدند. پس داریم:

$$MRR_{\text{QuranFastText}} = \frac{1}{2} \sum_{i=1}^2 \frac{1}{\text{Rank}_{\text{choice}}^i} = \frac{1}{2} [1 + 1] = 1.0$$

برای داده های نهج البلاغه هم داریم:

ورودی: الا و ان الشيطان قد جمع حربه و استجلب خيله و رجله و ان معى لبصيرتى ما لبست على نفسى و لا لبس على و ايم الله لافرن لهم حوزا انا ماتحه لا يصدرون عنه و لا يعودون اليه خروجی:

۱. الا و ان الشيطان قد جمع حربه و استجلب خيله و رجله و ان معى لبصيرتى ما لبست على نفسى و لا لبس على و ايم الله لافرن لهم حوزا انا ماتحه لا يصدرون عنه و لا يعودون اليه
۲. و ان على من الله جنة حصينه فاذا جاء يومى انفرجت عنى و اسلمتنى فحينئذ لا يطيش السهم و لا ييرا الكلم



۳. الحمد لله الذی لا توارى عنه سماء سماء و لا ارض ارضا منها و قد قال قائل انک على هذا الامر يا ابن ابی طالب لحريص فقلت بل انتم و الله لا حرص و ابد و انا اخص و اقرب و انما طلبت حقا لی و انتم تحولون بينی و بينه و تضربون وجهی دونه فلما قرعته بالحجه فی الملا الحاضرين هب کانه بهت لا یدری ما یجیبنی به اللهم انی استعديک على قریش و من اعانهم فانهم قطعوا رحمی و صغروا عظیم منزلتی و اجمعوا على منازعتی امرا هو لی ثم قالوا الا ان فی الحق ان تاخذه و فی الحق ان تترکه فخرجوا یجرون حرمه رسول الله-ص کما تجر الامة عند شرائها متوجهين بها الى البصره فحبسا نساءهما فی بیوتهما و ابرزا حبیس رسول الله-ص لهما و لغيرهما فی جيش ما منهم رجل الا و قد اعطانی الطاعه و سمح لی بالبیعه طائعا غیر مکره فقدموا على عاملی بها و خزان بیت مال المسلمین و غیرهم من اهلها فقتلوا طائفه صبرا و طائفه غدرا فوالله لو لم یصیبوا من المسلمین الا رجلا واحدا معتمدين لقتله بلا جرم جرّه لحل لی قتل ذلک الجيش کله اذ حضروه فلم ینکروا و لم یدفعوا عنه بلسان و لا بید دع ما انهم قد قتلوا من المسلمین مثل العده التي دخلوا بها علیهم

۴. من عبد الله على امیر المؤمنین الى القوم الذین غضبوا لله حین عصی فی ارضه و ذهب بحقه ف ضرب الجور سرادقه على البر و الفاجر و المقیم و الطاعن فلا معروف یستراح اليه و لا منکر یتناهی عنه اما بعد فقد بعثت الیکم عبدا من عباد الله لا ینام ایام الخوف و لا ینکل عن الاعداء ساعات الروح اشد على الفجار من حریق النار و هو مالک بن الحارث اخو مذحج فاسمعوا له و اطیعوا امره فیما طابق الحق فانه سیف من سیوف الله لا کلیل الظبه و لا نابی الضریبه فان امرکم ان تنفروا فانفروا و ان امرکم ان تقیموا فاقیموا فانه لا یقدم و لا یحجم و لا یؤخر و لا یقدم الا عن امری و قد اثرتکم به على نفسی لنصیحتة لکم و شده شکیمته على عدوکم

۵. ایها الناس غیر المغفول عنهم و التارکون الماخوذ منهم ما لی اراکم عن الله ذاهبین و الى غیره راغبین کانکم نعم اراح بها سائم الى مرعی و بی و مشرب دوی و انما هی کالمعلوفه للمدی لا تعرف ما ذا یراد بها اذا احسن اليها تحسب یومها دهرها و شبعها امرها و الله لو شئت ان اخبر کل رجل منکم بمخرجه و مولجه و جمیع شانه لفعلت و لكن اخاف ان تکفروا فی بر رسول الله-ص الا و انی مفضیه الى الخاصه ممن یؤمن ذلک منه و الذی بعثه بالحق و اصطفاه على الخلق ما انطق الا صادقا و قد عهد الى بذلک کله و بمهلک من یهلک و منجی من ینجو و مال هذا الامر و ما ابقى شیئا یمر على راسی الا افرغه فی اذنی و افضی به الى ایها الناس انی و الله ما احثکم على طاعه الا و اسبقکم اليها و لا انهاکم عن معصیه الا و اتناهی قبلکم عنها

۶. اتق الله الذی لا بد لک من لقائه و لا منتهی لک دونه و لا تقاتلن الا من قاتلک و سر البردین و غور بالناس و رفه فی السیر و لا تسر اول اللیل فان الله جعله سکنا و قدره مقاما لا ظعننا فارح فيه بدنک و روح ظهرک فاذا وقفت حین ینبطح السحر او حین ینفجر الفجر فسر على برکه الله فاذا لقیتم العدو فقف من اصحابک وسطا و لا تدن من

القوم دنو من یرید ان ینشب الحرب و لا تباعد عنهم تباعد من یهاب الباس حتی یاتیک امری و لا یحملنکم شنانهم علی قتالهم قبل دعائهم و الاعذار الیهم

۷. ما هی الا الکوفه اقبضها و ابسطها ان لم تکنی الا انت تهب اعاصیرک فقبحک الله و تمثل بقول الشاعر لعمر ابیک الخیر یا عمرو اننی || علی وضر من ذا الاناء قلیل ثم قال ع انبت بسرا قد اطلع الیمن و انی و الله لاظن ان هؤلاء القوم سیدالون منکم باجتماعهم علی باطلهم و تفرقکم عن حقکم و بمعصیتکم امامکم فی الحق و طاعتهم امامهم فی الباطل و بادائهم الامانه الی صاحبهم و خیانتکم و بصلاحهم فی بلادهم و فسادکم فلو ائتمنت احدکم علی قعب لخشیت ان یذهب بعلاقته اللهم انی قد مللتهم و ملونی و سئمتهم و سئموننی فابدلنی بهم خیرا منهم و ابدلهم بی شرا منی اللهم مٹ قلوبهم کما یماث الملح فی الماء اما و الله لوددت ان لی بکم الف فارس من بنی فراس بن غنم هنالك لو دعوت اتاک منهم || فوارس مثل ارمیه الحمیم ثم نزل ع من المنبر

۸. ۱۱

۹. و قد امرت علیکمما و علی من فی حیزکمما مالک بن الحارث الاشتر فاسمعا له و اطیعا و اجعلاه درعا و مجنا فانه ممن لا یخاف وهنه و لا سقطته و لا بطؤه عما الاسراع الیه احزم و لا اسراعه الی ما البطء عنه امثل ۱۰. و اعلم ان البصره مهبط ابلیس و مغرس الفتن فحادث اهلها بالاحسان الیهم و احلل عقده الخوف عن قلوبهم و قد بلغنی تنمرک لبنی تمیم و غلظتک علیهم و ان بنی تمیم لم یغب لهم نجم الا طلع لهم اخر و انهم لم یسبقوا بوغم فی جاهلیه و لا اسلام و ان لهم بنا رحما ماسه و قرابه خاصه نحن ماجورون علی صلتها و مازورون علی قطیعتها فاربع ابا العباس رحمک الله فیما جرى علی لسانک و یدک من خیر و شر فانا شریکان فی ذلک و کن عند صالح ظنی بک و لا یفیلن رایبی فیک و السلام

در این جا نفر اول جایگاه ۳ و نفر دوم جایگاه ۱ را حدس زد پس باز هم داریم:

$$MRR_{NahjFastText} = \frac{1}{2} \sum_{i=1}^2 \frac{1}{Rank_{choice}^i} = \frac{1}{2} \left[ \frac{1}{3} + 1 \right] = 0.665$$

برای صحیفه سجادیه نیز همین عملیات صورت گرفت:

ورودی: حمدا تقر به عیوننا اذا برقت الابصار، و تبيض به وجوهنا اذا اسودت الابصار.

خروجی:

۱. حمدا تقر به عیوننا اذا برقت الابصار، و تبيض به وجوهنا اذا اسودت الابصار.

۲. حمدا نعمر به فیمن حمده من خلقه، و نسبق به من سبق الی رضاه و عفوه.

## یافتن عبارات مشابه در متون دینی با استفاده از جاگذاری کلمات

تکلیف سری سوم

۳. یا من تحمد الی عبادہ بالاحسان و الفضل، و غمرهم بالمن و الطول، ما افشی فینا نعمتک، و اسبغ علینا منتک، و اخصنا ببرک!
۴. اللهم و قد اشرف علی خفایا الاعمال علمک، و انکشف کل مستور دون خبرک، و لا تنطوی عنک دقائق الامور، و لا تعذب عنک غیبات السرائر
۵. و هب لی التطهیر من دنس العصیان، و اذهب عنی درن الخطایا، و سربلنی بسربال عافیتک، و ردنی رداء معافاتک، و جللنی سوابغ نعماتک، و ظاهر لدی فضلک و طولک
۶. و انا ابرا الیک من ان استکبر، و اعوذ بک من ان اصر، و استغفرک لما قصرت فیه، و استعین بک علی ما عجزت عنه.
۷. و الذین عرفتهم مثاقیل المیاه، و کیل ما تحویه لواعج الامطار و عوالجها
۸. و اذا هممنا بهمین یرضیک احدهما عنا، و یسخطک الاخر علینا، فمل بنا الی ما یرضیک عنا، و اوھن قوتنا عما یسخطک علینا
۹. لا یخیب منک الاملون، و لا ییاس من عطائک المتعرضون، و لا یشقی بنقمتک المستغفرون.
۱۰. و اجزل لنا فیه من الحسنات، و اخلنا فیه من السيئات، و املا لنا ما بین طرفیه حمدا و شکرا و اجرا و ذخرا و فضلا و احسانا.

نفر اول جایگاه ۱ و نفر دوم نیز جایگاه ۲ را حدس زد که داریم:

$$MRR_{Sahifa_{FastText}} = \frac{1}{2} \sum_{i=1}^2 \frac{1}{Rank_{choice}^i} = \frac{1}{2} \left[ \frac{1}{2} + 1 \right] = 0.75$$

و در کل هم شاخص MRR را به صورت زیر برای مدل FastText داریم:

$$MRR_{TF-IDF} = \frac{0.75 + 0.665 + 1}{3} = 80\%$$

## FINE-TUNED BERT – استفاده از شبکه‌های

### آموزش دادن مدل برت

برای لود کردن مدل برت از کتابخانه transformers و از یکی از مدل آموزش دیده روی زبان عربی (البته بیس ورژن آن) استفاده شده است.

```
from transformers import AutoTokenizer, AutoModel
import torch
device = 'cuda' if torch.cuda.is_available() else 'cpu'

# Mini: asafaya/bert-mini-arabic
# Medium: asafaya/bert-medium-arabic
# Base: asafaya/bert-base-arabic
# Large: asafaya/bert-large-arabic
# https://www.youtube.com/watch?v=jVPd7lEvjtg

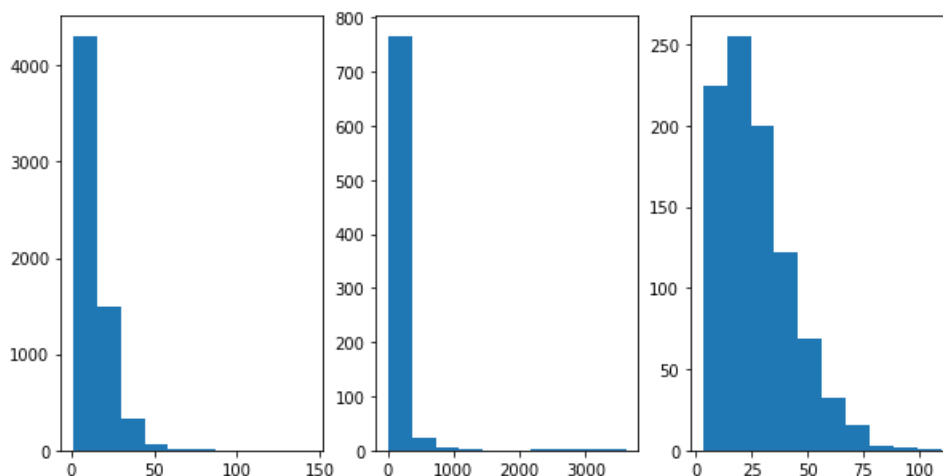
tokenizer = AutoTokenizer.from_pretrained("asafaya/bert-base-arabic")
model = AutoModel.from_pretrained("asafaya/bert-base-arabic")

model=model.to(device)

model.eval()

device
```

ابتدا یک پلات کلی برای هیستوگرام تعداد توکن‌ها در هر جمله از قران و صحیفه سجادیه و نهج البلاغه کشیده شد. به نظر می‌آید مناسب‌ترین عدد ۱۲۸ باشد.



# یافتن عبارات مشابه در متون دینی با استفاده از جاگذاری کلمات

تکلیف سری سوم

بدین ترتیب با استفاده از توکنایزر برت، جملاتی که طول کمتر از ۱۲۸ دارند پد می‌شوند و آن‌ها که طول بیشتر از ۱۲۸ دارند بریده می‌شوند که در کد زیر نشان داده شده است:

```
def tokenize(sentence_list):
    return tokenizer.batch_encode_plus([sentence_list,
                                        add_special_tokens=True,
                                        max_length=128,
                                        padding='max_length',
                                        return_attention_mask=True,
                                        truncation=True,
                                        return_tensors='pt'])
```

تابع بعدی که در اینجا مهم است، تابع `get_mean_embedding` است. این تابع جملات خام و مدل برت را به عنوان ورودی می‌گیرد و یک امبدینگ به ازای هر جمله خروجی می‌دهد. برای بدست آوردن خروجی مدل برت از `last_hidden_state` استفاده شده است سپس روی کلماتی که واقعا در جمله بوده اند میانگین گرفته شده است. این مسئله به شدت اهمیت دارد زیرا در غیراینصورت بردارهای میانگین تقریبا شبیه به هم خواهند شد و اطلاعات خاصی را نشان نخواهند داد.

```
import numpy as np
import torch
from torch.utils.data import DataLoader, TensorDataset

def get_mean_embedding(model, sentences, batch_size=128):
    encoding = tokenize(sentences)

    dataset = TensorDataset(encoding.input_ids, encoding.attention_mask)
    dataloader = DataLoader(dataset, batch_size=batch_size,
                            drop_last=False, shuffle=False)

    predictions = []
    with torch.no_grad():
        for (input_ids, attention_mask) in dataloader:
            prediction = model(input_ids.to(device),
                               attention_mask.to(device))
            predictions.append(prediction.last_hidden_state)

    predictions = torch.cat(predictions, dim=0)

    attention_mask = encoding.attention_mask.unsqueeze(-1) \
        .expand(predictions.shape).float()
    # for filtering out unnecessary ones
    mask_embeddings = predictions.cpu() * attention_mask

    summed = torch.sum(mask_embeddings, dim=1)
    count = torch.clamp([attention_mask.sum(dim=1), min=1e-9])
    mean_embedding = summed / count

    return mean_embedding
```

سپس بردار میانگین جملات مختلف با یکدیگر مقایسه شده و شبیه‌ترین موارد برگردانده می‌شود. در اینجا از کتابخانه `sklearn` استفاده شده است.

## یافتن عبارات مشابه در متون دینی با استفاده از جاگذاری کلمات

تکلیف سری سوم

```
from sklearn.metrics.pairwise import cosine_similarity

def get_most_similar(mean_embedding, compare_verse_no):
    mean_embedding = mean_embedding.numpy()
    similarity = cosine_similarity(
        [mean_embedding[compare_verse_no]],
        mean_embedding
    ).squeeze()

    most_similar = similarity.argsort()[-10:][::-1]
    return most_similar
```

تا اینجا می‌توان یک سری نتایج بدست آورد ولی برای این که مدل را به خوبی فاین تیون کنیم ابتدا از تسک masked language modelling استفاده کردیم. این تسک بدین صورت است که ۱۵ درصد کلمات در هر جمله mask می‌شوند سپس مدل باید روی این کلمات آموزش ببیند. در واقع بدین صورت دارد کلمات تخصصی و پرکاربرد هر متن را یاد می‌گیرد که میتواند دقت بالاتری بدهد.

```
from transformers import BertTokenizer, BertForMaskedLM
import torch
from transformers import AdamW
device = 'cuda' if torch.cuda.is_available() else 'cpu'

# Mini: asafaya/bert-mini-arabic
# Medium: asafaya/bert-medium-arabic
# Base: asafaya/bert-base-arabic
# Large: asafaya/bert-large-arabic

tokenizer = BertTokenizer.from_pretrained("asafaya/bert-base-arabic")
model = BertForMaskedLM.from_pretrained("asafaya/bert-base-arabic")

model = model.to(device)
model.train()

optim = AdamW(model.parameters(), lr=5e-5)
```

سپس مدل با یک حلقه معمولی‌ترین می‌شود. تفاوت مدل ما با مدل اصلی برت (که نتایج بهتری هم می‌گیرد) این است که در هر ۵ اپاک داده‌های مسک شده مجدداً ساخته می‌شوند تا مدل روی مقدار داده بیش‌برازش نشود و بتواند قسمت‌های مختلف متن را یاد بگیرد. در نهایت در بخش بعد ورودی مدل را مشخص کرده تا در ادامه از هر کتاب، ۱۵ آیه/جمله مشابه توسط مدل فاین تیون شده bert مشخص شود.

برای محاسبه MRR به جهت مقایسه این روش شباهت‌یابی و دو روش بعدی، ترتیب ۱۵ آیه/جمله مشابه را بهم ریخته و به دو نفر اعضای تیم که این خروجی‌ها را ندیده‌اند، داده می‌شود تا به صورت دستی مشابه‌ترین آیه/جمله مدنظرشان را

# یافتن عبارات مشابه در متون دینی با استفاده از جاگذاری کلمات

تکلیف سری سوم

اعلام کنند. ورودی که برای روش سوم یعنی Fine-tuned BERT آزمایش شد، یکی از آیه‌های قرآن یعنی آیه ۲۴ سوره بقره: «فان لم تفعلوا ولن تفعلوا فاتقوا النار التي وقودها الناس والحجاره ۚ اعدت للكافرين» بود. توجه شود که چون یکی از آیه‌های قرآن را به عنوان ورودی مشخص کردیم، طبیعی است که برترین و مشابه‌ترین آیه قرآنی با این آیه، خودش باشد لذا برای قرآن، ما ۱۴ آیه مشابه بعدی را برای محاسبه دستی استفاده کردیم. ۱۴ آیه برتر در شباهت به ترتیب عبارتند از:

۱- یا ایها الذین امنوا قوا انفسکم واهلیکم نارا وقودها الناس والحجاره علیها ملائکه غلاظ شداد لا یعصون الله ما امرهم ویفعلون ما یأمرون

۲- یا قوم ادخلوا الارض المقدسه التي کتب الله لکم ولا ترتدوا علی ادبارکم فتنقلبوا خاسرین

۳- لن تنالوا البر حتی تنفقوا مما تحبون ۚ وما تنفقوا من شیء فان الله به علیم

۴- وما کان صلاتهم عند البیت الا مکاء وتصدیه ۚ فذوقوا العذاب بما کنتم تکفرون

۵- اشفقتم ان تقدموا بین یدی نجواکم صدقات ۚ فاذا لم تفعلوا وتاب الله علیکم فاقیموا الصلاه واتوا الزکاه واطیعوا الله ورسوله ۚ والله خبیر بما تعملون

۶- وجعلوا لله اندادا لیضلوا عن سبيله ۚ قل تمتعوا فان مصیرکم الی النار

۷- الذین کفروا لهم عذاب شدید ۚ والذین امنوا وعملوا الصالحات لهم مغفره واجر کبیر

۸- والذین کفروا بعضهم اولیاء بعض ۚ الا تفعلوه تکن فتنه فی الارض وفساد کبیر

...

۱۳- قل للذین کفروا ستغلبون وتحشرون الی جهنم ۚ وبئس المهاد

۱۴- واتقوا فتنه لا تصیبن الذین ظلموا منکم خاصه ۚ واعلموا ان الله شدید العقاب

نفر اول آیه شماره ۱ و نفر دوم آیه شماره ۲ را به عنوان مشابه‌ترین آیه برگزیدند. بنابراین داریم:

$$MRR_{QuranBERT} = \frac{1}{2} \sum_{i=1}^2 \frac{1}{Rank_{choice}^i} = \frac{1}{2} \left[ 1 + \frac{1}{2} \right] = 0.75$$

به همین ترتیب آیه/جمله‌های برتر در نهج‌البلاغه عبارت بودند از:

۱- لا تلقین طلحه فانک ان تلقه تجده کالثور عاقصا قرنه یرکب الصعب و یقول هو الذلول و لكن الق الزبیر فانه الین عریکه  
فقل له یقول لک ابن خالک عرفتنی بالحجاز و انکرتنی بالعراق فما عدا مما بدا

۲- ما کنت تصنع بسعه هذه الدار فی الدنيا و انت الیها فی الاخره کنت احوج و بلی ان شئت بلغت بها الاخره تقری فیها  
الضیف و تصل فیها الرحم و تطلع منها الحقوق مطالعها فاذا انت قد بلغت بها الاخره فقل له العلاء یا امیر المؤمنین اشکو  
الیک اخی عاصم بن زیاد قال و ما له قال لبس العباءه و تخلی عن الدنيا قال علی به فلما جاء قال یا عدی نفسه لقد  
استهام بک الخبیث ا ما رحمت اهلک و ولدک ا ترى الله احل لک الطیبات و هو یکره ان تاخذها انت اهون علی الله من  
ذلک قال یا امیر المؤمنین هذا انت فی خشونه ملبسک و جشوبه ما کلک قال و یحک انی لست کانت ان الله تعالی فرض  
علی ائمه العدل ان یقدروا انفسهم بضعفه الناس کیلا یتبیغ بالفقیر فقر

۳- و ساله رجل ان یعرفه الایمان فقال ع اذا کان الغد فاتنی حتی اخبرک علی اسماع الناس فان نسیت مقالتی حفظها  
علیک غیرک فان الکلام کالشارده ینقفها هذا و یخطئها هذا

...

۱۴- ایها الناس من عرف من اخیه وثیقہ دین و سداد طریق فلا یسمعن فیہ اقاولیل الرجال اما انه قد یرمی الرامی و  
تخطئ السهام و یحیل الکلام و باطل ذلک یبور و الله سمیع و شهید اما انه لیس بین الحق و الباطل الا اربع اصابع الباطل  
ان تقول سمعت و الحق ان تقول رایت

۱۵- اما بعد فان صلاح ابیک غرنی منک و ظننت انک تتبع هدیه و تسلك سبيله فاذا انت فیما رقی الی عنک لا تدع  
لهواک انقیادا و لا تبقى لاخرتک عتادا تعمر دنیاک بخراب اخرتک و تصل عشیرتک بقطیعه دینک و لئن کان ما بلغنی  
عنک حقا لجمال اهلک و شسع نعلک خیر منک و من کان بصفتک فلیس باهل ان یسد به ثغر او ینفذ به امر او یعلی له  
قدر او یشرک فی امانه او یؤمن علی جبايه فاقبل الی حین یصل الیک کتابی هذا ان شاء الله

نفر اول آیه شماره ۳ و نفر دوم آیه شماره ۱ را به عنوان مشابه‌ترین آیه برگزیدند. بنابراین داریم:

$$MRR_{NahjBERT} = \frac{1}{2} \sum_{i=1}^2 \frac{1}{Rank_{choice}^i} = \frac{1}{2} \left[ 1 + \frac{1}{3} \right] = 0.665$$



# یافتن عبارات مشابه در متون دینی با استفاده از جاگذاری کلمات

تکلیف سری سوم

به همین ترتیب آیه/جمله‌های برتر در صحیفه سجادیه عبارت بودند از:

۱- فختم بنا علی جمیع من ذرا، و جعلنا شهداء علی من جحد، و کثرنا بمنه علی من قل.

۲- هذا مقام من اعترف بسبوغ النعم، و قابلها بالتقصير، و شهد علی نفسه بالتضييع.

۳- اللهم انک انزلته علی نبیک محمد- صلی الله علیه و اله- مجملا، و الهمته علم عجائب مکملا، و ورثتنا علمه مفسرا، و فضلنا علی من جهل علمه، و قویتنا علیه لترفعنا فوق من لم یطق حمله.

...

۱۴- الحمد لله الذی هدانا لحمده، و جعلنا من اهله لنکون لاحسانه من الشاکرین، و لیجزینا علی ذلک جزاء المحسنین.

۱۵- السلام علیک ما کان اطولک علی المجرمین، و اهییک فی صدور المؤمنین!

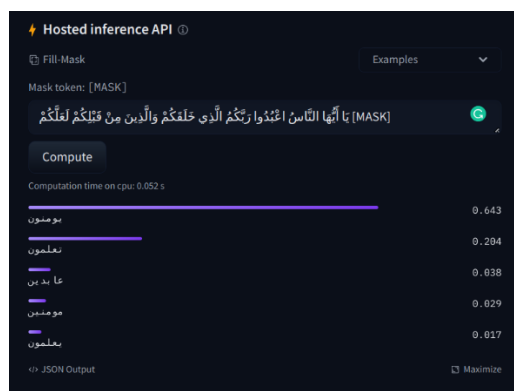
نفر اول آیه شماره ۲ و نفر دوم آیه شماره ۲ را به عنوان مشابه‌ترین آیه برگزیدند. بنابراین داریم:

$$MRR_{SahifaBERT} = \frac{1}{2} \sum_{i=1}^2 \frac{1}{Rank_{choice}^i} = \frac{1}{2} \left[ \frac{1}{2} + \frac{1}{2} \right] = 0.5$$

بنابراین اگر روش Fine-tuned Bert را بخواهیم با یک MRR نشان دهیم، با میانگین گیری از این سه عدد بدست آمده داریم:

$$MRR_{BERT} = \frac{0.75 + 0.665 + 0.5}{3} = 63.8\%$$

البته مدل ما در سایت Huggingface نیز آپلود شده است و بصورت آنلاین از طریق [این لینک](#) قابل استفاده است پیشنهاد می‌کنیم حتما به آن سری بزنید.



## یافتن عبارات مشابه در متون دینی با استفاده از جاگذاری کلمات

تکلیف سری سوم