



به نام خدا



دانشگاه تهران

دانشکده مهندسی برق و کامپیوتر

پردازش زبان‌های طبیعی

تمرین شماره 5

نام و نام خانوادگی	محمد جواد رنجبر
شماره دانشجویی	۸۱۰۱۰۱۱۷۳
تاریخ ارسال گزارش	۱۴۰۲/۰۳/۱۵

فهرست گزارش

سوال ۱	۵
پیش پردازش دادگان با استفاده از preprocess-fairseq	۵
آموزش مدل با استفاده از train-fairseq	۵
استفاده از generate-fairseq برای ترجمه دادگان ارزیابی	۵
استفاده از Tensorboard برای بررسی فرایند آموزش	۷
آموزش مدل BPE برای دادگان انگلیسی و فارسی	۸
پردازش دادگان با مدل های BPE	۸
آموزش مدل با استفاده از train-fairseq	۸
استفاده از generate-fairseq برای ترجمه دادگان ارزیابی	۸
استفاده از Tensorboard برای بررسی فرایند آموزش	۱۰
سوال ۲	۱۱
دانلود Bert Model و Bert Tokenizer	۱۱
پردازش دادگان AFEC با استفاده از Bert Tokenizer	۱۱
پیش پردازش دادگان با استفاده از preprocess-fairseq	۱۱
ذخیره وزن های اولیه embedding شبکه Bert با فرمت مناسب	۱۱
آموزش مدل با استفاده از fairseq-train	۱۱
Freeze شده	۱۱
Freeze نشده	۱۲
استفاده از generate-fairseq برای ترجمه دادگان ارزیابی	۱۲
Freeze شده	۱۲
Freeze نشده	۱۳

۱۴..... استفاده از Tensorboard برای بررسی فرایند آموزش

۱۴..... Freeze شده

۱۵..... Freeze نشده

فهرست شکل ها

- شکل ۱ loss بدون استفاده از bpe ۷
- شکل ۲ BLEU بدون BPE ۷
- شکل ۳ loss مدل با توکنایزشن BPE ۱۰
- شکل ۴ BLEU با توکنایزشن BLEU ۱۰
- شکل ۵ loss مدل freeze شده و استفاده از Bert embedding ۱۴
- شکل ۶ Bleu مدل freeze شده و استفاده از Bert embedding ۱۵
- شکل ۷ loss مدل freeze نشده و استفاده از Bert embedding ۱۵
- شکل ۸ Bleu مدل freeze نشده و استفاده از Bert embedding ۱۶

تمام فایل‌های خواسته شده شامل CSVها در فولدر مربوط قرار داده شده است.

سوال ۱

پیش پردازش دادگان با استفاده از preprocess-fairseq

در این مرحله این کارها توسط fairseq انجام می‌شود:

توکن‌بندی: Fairseq عمل توکن‌بندی را انجام می‌دهد که در آن متن به واحدهای جداگانه‌ای مانند

کلمات، زیرکلمات یا حروف تقسیم می‌شود.

ایجاد واژگان: Fairseq واژگان بر اساس داده‌های توکن‌بندی شده ایجاد می‌کند. در این تمرین ۴۰۰۰۰

واژه تولید می‌کنیم.

عددسازی: بعد از توکن‌بندی، Fairseq یک نمایش عددی به هر توکن اختصاص می‌دهد.

آموزش مدل با استفاده از train-fairseq

پارامترهای داده شده در سوال را تنظیم می‌کنیم:

معماری یک لایه encoder و decoder و یک لایه lstm با مکانیزم attention

تابع هزینه label smooth cross entropy به مقدار label smoothing برابر ۰.۲

نرخ یادگیری تطبیقی inverse sqrt با مقدار اولیه ۰.۰۰۲۵

بهینه‌ساز adam با پارامترهای $\beta_1=0.9$ و $\beta_2=0.98$

Dropout با مقدار ۰.۲۵

برای ۶ epoch آموزش گذاشتیم.

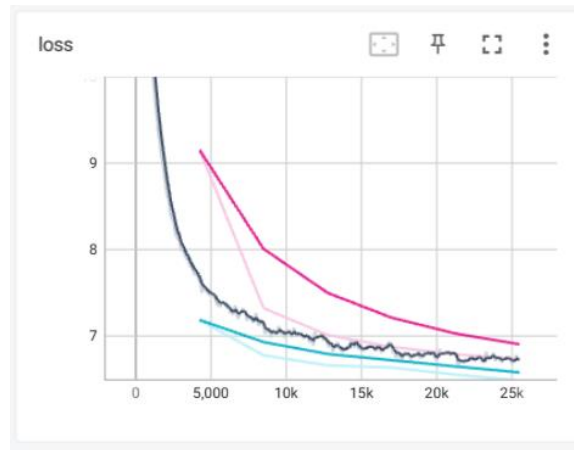
استفاده از generate-fairseq برای ترجمه دادگان ارزیابی

تعدادی نمونه از این ترجمه‌ها به صورت زیر می‌باشد:

S-7043	it is not bad either ,	
T-7043	این هم بد نیست ,	
H-7043	-1.0841904878616333	این بد نیست ,
D-7043	-1.0841904878616333	این بد نیست ,
P-7043	-2.7383 -1.1646 -0.4039 -0.7685 -0.3457	
S-7021	yes , Friday the fifth .	
T-7021	بله , جمعه پنجم .	
H-7021	-0.5548604726791382	بله , جمعه پنجم .
D-7021	-0.5548604726791382	بله , جمعه پنجم .
P-7021	-0.1186 -0.1450 -1.9711 -0.2583 -0.4583 -0.3779	
S-6924	yes , that is okay .	
T-6924	بله , خوب است .	
H-6924	-0.5492016077041626	بله , خوب است .
D-6924	-0.5492016077041626	بله , خوب است .
P-6924	-0.1230 -0.1799 -2.3391 -0.1724 -0.2090 -0.2718	
S-6920	which Friday do you mean ?	
T-6920	منظورتان کدام جمعه است	
H-6920	-1.489291787147522	? که جمعه به معنای آن است
D-6920	-1.489291787147522	? که جمعه به معنای آن است
P-6920	-2.0477 -0.1897 -3.9687 -1.8533 -2.1071 -0.8938 -0.6505 -0.2035	
S-6911	yes , at what time ?	
T-6911	? بله , چه ساعتی	

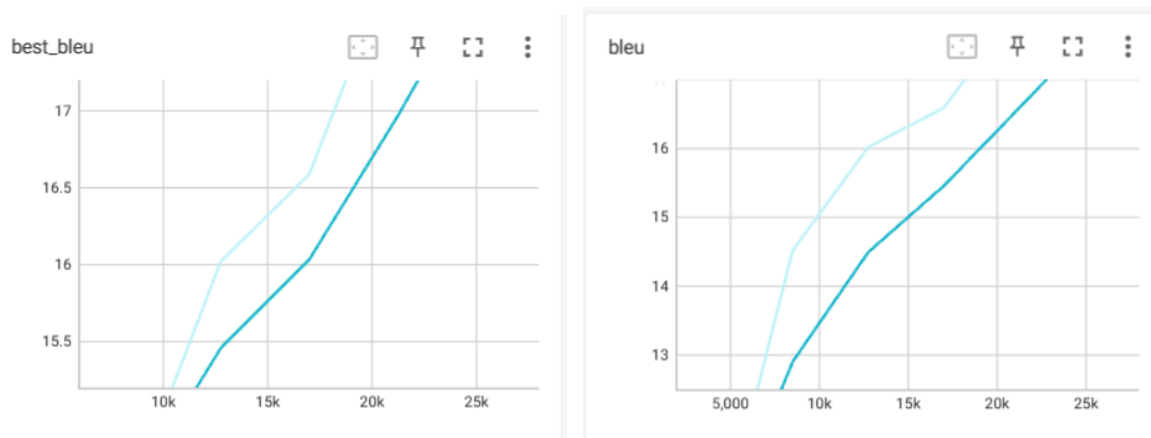
استفاده از Tensorboard برای بررسی فرایند آموزش

loss به شکل زیر است:



شکل ۱ loss بدون استفاده از bpe

نمودار BLEU برای دادگان ارزیابی به شکل‌های زیر می‌باشد:



شکل ۲ BLEU بدون BPE

همچنین مقدار نهایی این تابع برای داده‌های validation به شکل زیر است.

Generate test with beam=5: BLEU4 = 24.73, 58.5/31.7/18.3/11.0 (BP=1.000, ratio=1.093, syslen=99070, reflen=90650)

بر اساس جدول زیر می‌توان راجع به گفت مدل تقریباً مناسب کار می‌کند.

BLEU Score

< 10
10 - 19
20 - 29

Interpretation

Almost useless
Hard to get the gist
The gist is clear, but has significant grammatical errors

BLEU Score	Interpretation
30 - 40	Understandable to good translations
40 - 50	High quality translations
50 - 60	Very high quality, adequate, and fluent translations
> 60	Quality often better than human

آموزش مدل BPE برای دادگان انگلیسی و فارسی

ابتدا با استفاده از SentencePieceTrainer کلمات را با روش BPE کلات را توکنایز می‌کنیم.

پردازش دادگان با مدل های BPE

حال همانند بخش قبل داده‌ها را preprocess می‌کنیم که کارهایی که انجام می‌شد، ذکر شده است.

آموزش مدل با استفاده از train-fairseq

پارامترهای داده شده در سوال را تنظیم می‌کنیم:

معماری یک لایه encoder و decoder و یک لایه lstm با مکانیزم attention

تابع هزینه label smooth cross entropy به مقدار label smoothing برابر ۰.۲

نرخ یادگیری تطبیقی inverse sqrt با مقدار اولیه 0.0025

بهینه‌ساز adam با پارامترهای $\beta_1=0.9$ و $\beta_2=0.98$

Dropout با مقدار ۰.۲۵

برای ۶ epoch آموزش گذاشتیم.

استفاده از generate-fairseq برای ترجمه دادگان ارزیابی

نمونه‌ای از جملات ترجمه شده به صورت زیر می‌باشد.

S-3414 __that __would __be __fine.__

T-3414 __ خوب __ است __.

H-3414 __ که __ خوب __ است __. -۰.۷۵۹۴۴۹۰۰۵۱۲۶۹۵۳۱-

D-3414 __ که __ خوب __ است __. -۰.۷۵۹۴۴۹۰۰۵۱۲۶۹۵۳۱-

- P-3414 ۰.۲۵۹۴- ۰.۴۱۲۳- ۱.۰۹۰۴- ۰.۳۵۸۶- ۰.۵۵۱۷- ۰.۲۸۴۴- ۲.۸۵۸۹- ۰.۲۵۹۹-
- S-3406 _do _you _know _that?_
- T-3406 >> _unk?_ << _آنرا _
- H-3406 ۰.۸۰۳۱۸۴۷۴۷۶۹۵۹۲۲۹- _آیا _ می _ دانید _?
- D-3406 ۰.۸۰۳۱۸۴۷۴۷۶۹۵۹۲۲۹- _آیا _ می _ دانید _?
- P-3406 ۰.۲۲۰۱- ۱.۲۷۶۲- ۰.۵۴۶۹- ۰.۳۴۱۴- ۰.۹۳۴۳- ۰.۳۳۹۵- ۲.۴۲۴۷- ۰.۳۴۲۳-
- S-2996 _at _two _P M._
- T-2996 _ دو _ بعدازظهر _.
- H-2996 ۰.۸۰۱۹۴۶۹۳۸۰۳۷۸۷۲۳- _ در _ دو _ بعدازظهر _.
- D-2996 ۰.۸۰۱۹۴۶۹۳۸۰۳۷۸۷۲۳- _ در _ دو _ بعدازظهر _.
- P-2996 ۰.۳۰۰۳- ۰.۷۶۵۱- ۲.۴۶۸۴- ۰.۳۶۳۶- ۰.۳۱۵۰- ۰.۳۰۳۴- ۱.۶۳۰۵- ۰.۲۶۹۳-
- S-2710 _that _would _be _great._
- T-2710 _ عالی _ است _.
- H-2710 ۰.۹۵۰۵۳۹۹۴۶۵۵۶۰۹۱۳- _ این _ عالی _ است _.
- D-2710 ۰.۹۵۰۵۳۹۹۴۶۵۵۶۰۹۱۳- _ این _ عالی _ است _.
- P-2710 ۰.۲۳۸۷- ۰.۴۱۶۹- ۱.۲۸۴۲- ۰.۳۳۷۲- ۱.۱۷۸۹- ۰.۳۱۱۰- ۳.۵۷۵۱- ۰.۲۶۲۴-
- S-2660 _you _will _find _that._
- T-2660 _ پیدا _ میکنید _.
- H-2660 ۰.۸۰۹۲۳۱۹۹۶۵۳۶۲۵۴۹- _ شما _ پیدا _ خواهید _ کرد _.
- D-2660 ۰.۸۰۹۲۳۱۹۹۶۵۳۶۲۵۴۹- _ شما _ پیدا _ خواهید _ کرد _.
- P-2660 - ۰.۱۷۶۵- ۰.۲۸۹۵- ۰.۹۹۱۸- ۰.۲۴۶۵- ۳.۷۴۷۴- ۰.۲۵۵۴- ۱.۳۸۴۱- ۰.۲۹۲۴-
۰.۲۴۸۳- ۰.۴۶۰۵
- S-2676 _the _end _of _May._

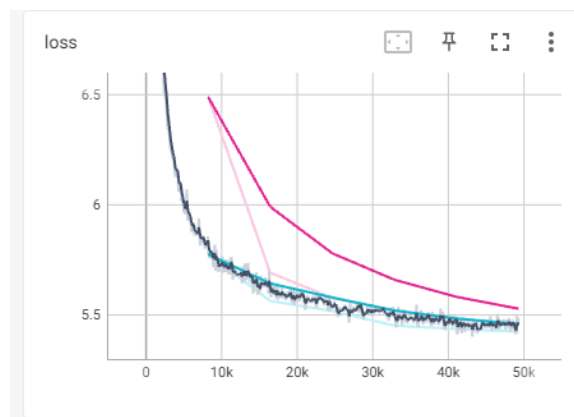
T-2676 — اواخر — می — .

H-2676 — پایان — ماه — می — .۰.۵۳۱۲۶۱۸۶۱۳۲۴۳۱۰۳-

D-2676 — پایان — ماه — می — .۰.۵۳۱۲۶۱۸۶۱۳۲۴۳۱۰۳-

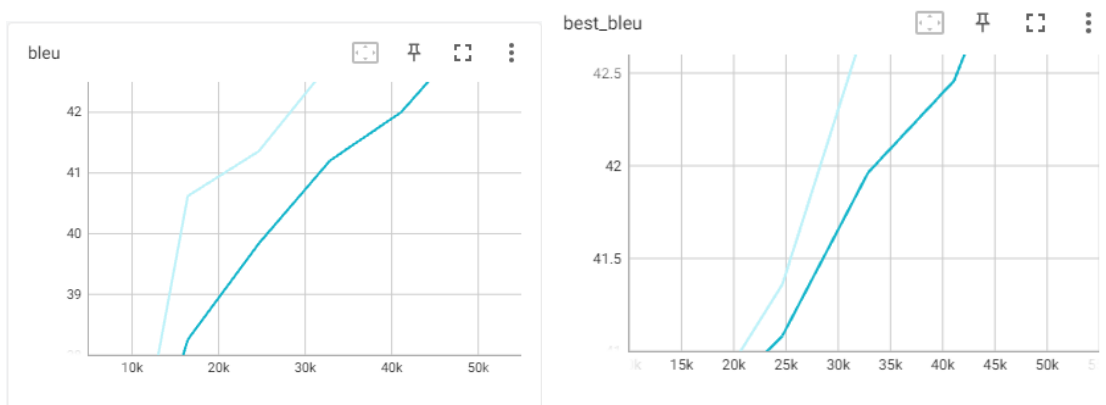
استفاده از Tensorboard برای بررسی فرایند آموزش

Loss به شکل زیر می باشد:



شکل ۳ loss مدل با توکنایزشن BPE

همچنین نمودار BLEU به شکل زیر خواهد بود:



شکل ۴ BLEU با توکنایزشن BLEU

همچنین مقدار نهایی BLEU برای داده های ارزیابی به صورت زیر می باشد:

Generate test with beam=5: BLEU4 = 42.91, 69.4/49.7/37.8/26.0 (BP=1.000, ratio=1.093, syslen=186038, reflen=170167)

این مقدار نشان می دهد که مدل به شدت خوب عمل می کند و BPE روش بهتری بود.

سوال ۲

دانلود Bert Model و Bert Tokenizer

ابتدا با استفاده از کتابخانه transformers مدل Bert multilingual را لود می‌کنیم.

پردازش دادگان AFEC با استفاده از Bert Tokenizer

حال داده‌ها را به ترتیب با استفاده از این مدل tokenize می‌کنیم و در پوشه‌ای جدید با نام Tokenized_data ذخیره می‌کنیم.

پیش پردازش دادگان با استفاده از preprocess-fairseq

مانند سوال قبل نیز در این مرحله مراحل پیش‌پردازش توسط fairseq انجام می‌شود.

ذخیره وزن‌های اولیه embedding شبکه Bert با فرمت مناسب

از آنجا که تعداد کلمات مدل bert خیلی زیاد بود و colab با این تعداد embedding را نمی‌توانست در حافظه نگه دارد، فقط ۵۰۰۰ تا از کلمه‌ها را نگه می‌داریم. و بر اساس فرمت خواسته شده در سوال این کلمات را ذخیره می‌کنیم.

آموزش مدل با استفاده از fairseq-train

ابتدا با استفاده از دستورات زیر مدل را برای آموزش با embedding استخراج شده از Bert را لود کرده و بر اساس آن‌ها آموزش می‌دهیم. در این مرحله یک بار مدل را freeze کردن و یک بار بدون این کار آموزش می‌دهیم،

Freeze شده

در این نحوه با استفاده از دستورات زیر بخش‌هایی از مدل را فریز می‌کنیم.

```
--encoder-freeze-embed \  
--decoder-freeze-embed \  

```

Freeze نشده

استفاده از generate-fairseq برای ترجمه دادگان ارزیابی

تعدادی از نتایج به صورت زیر می باشد.

Freeze شده

S-6392	that is oka ##y.
T-6392	خوب است .
H-6392	آن بسیار خوب است .
D-6392	آن بسیار خوب است .
P-6392	۵.۱۸۹۷- ۳.۱۵۱۳- ۲.۴۶۵۲- ۰.۱۱۹۱- ۰.۲۰۶۵- ۰.۳۱۳۱-
S-6370	that is oka ##y.
T-6370	خوب است .
H-6370	آن بسیار خوب است .
D-6370	آن بسیار خوب است .
P-6370	۵.۱۸۹۷- ۳.۱۵۱۳- ۲.۴۶۵۲- ۰.۱۱۹۱- ۰.۲۰۶۵- ۰.۳۱۳۱-
S-6321	is that oka ##y?
T-6321	خوب است ؟
H-6321	این م ##سا ##له است ؟
D-6321	این م ##سا ##له است ؟
P-6321	۲.۷۶۴۸- ۶.۴۲۳۴- ۲.۴۰۸۳- ۰.۱۷۸۷- ۰.۷۱۹۷- ۰.۲۴۱۲- ۰.۲۰۴۴-
S-5998	that is oka ##y.
T-5998	خوب است .
H-5998	آن بسیار خوب است .
D-5998	آن بسیار خوب است .

P-5998 ۰.۳۱۳۱- ۰.۲۰۶۵- ۰.۱۱۹۱- ۲.۴۶۵۲- ۳.۱۵۱۳- ۵.۱۸۹۷-

S-8503 make a suggest ##ion.

T-8503 پیشنهاد ب ##ده .

Freeze نشده

P-7234 - ۰.۳۸۷۹- ۱.۶۹۸۷- ۰.۰۸۸۰- ۰.۹۴۶۲- ۱.۸۶۲۰- ۰.۳۸۲۴- ۰.۲۲۳۳- ۰.۱۷۵۷-
۰.۲۷۵۷

S-7012 in the afternoon it is still possible.

T-7012 بعد از ظهر میشود .

H-7012 ۰.۸۳۸۹۸۸۹۵۹۷۸۹۲۷۶۱- در بعد ## از ##ظهر هنوز امکان ##پ ##ذیر است .

D-7012 ۰.۸۳۸۹۸۸۹۵۹۷۸۹۲۷۶۱- در بعد ## از ##ظهر هنوز امکان ##پ ##ذیر است .

P-7012 - ۰.۱۱۹۵- ۰.۴۴۱۹- ۰.۹۶۹۷- ۱.۸۹۳۹- ۰.۳۶۶۸- ۱.۵۷۹۲- ۰.۵۲۷۸- ۲.۳۰۹۲-
۰.۲۷۸۰- ۰.۳۸۱۶- ۰.۳۶۱۴

S-6559 ye ##s , that would be possible.

T-6559 بل ##ه , میشود .

H-6559 ۰.۶۷۱۱۰۷۳۵۱۷۷۹۹۳۷۷- بل ##ه , امکان ##پ ##ذیر است .

D-6559 ۰.۶۷۱۱۰۷۳۵۱۷۷۹۹۳۷۷- بل ##ه , امکان ##پ ##ذیر است .

P-6559 - ۰.۳۸۷۹- ۱.۶۹۸۷- ۰.۰۸۸۰- ۰.۹۴۶۲- ۱.۸۶۲۰- ۰.۳۸۲۴- ۰.۲۲۳۳- ۰.۱۷۵۷-
۰.۲۷۵۷

S-6229 that is a very good suggest ##ion.

T-6229 پیشنهاد بسیار خوبی است .

H-6229 ۱.۰۶۲۴۰۹۲۸۱۷۳۰۶۵۱۹- این یک پیشنهاد بسیار خوب است .

D-6229 ۱.۰۶۲۴۰۹۲۸۱۷۳۰۶۵۱۹- این یک پیشنهاد بسیار خوب است .

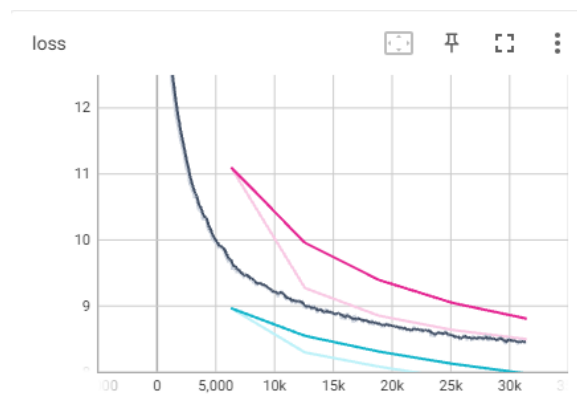
P-6229 ۰.۳۱۷۹- ۰.۴۹۶۹- ۰.۵۷۹۹- ۰.۵۶۴۷- ۱.۲۷۷۳- ۱.۱۰۳۹- ۱.۶۳۶۸- ۲.۵۲۱۹-

S-2330	ten o ' clock in the morning.
T-2330	ساعت ده ص ###یح .
H-2330	ساعت ده ص ###یح . -۰.۴۷۰۹۹۸۰۴۸۷۸۲۳۴۸۶۳
D-2330	ساعت ده ص ###یح . -۰.۴۷۰۹۹۸۰۴۸۷۸۲۳۴۸۶۳
P-2330	-۰.۴۶۴۰- ۱.۸۰۱۴- ۰.۰۵۴۹- ۰.۰۶۰۴- ۰.۳۸۰۲- ۰.۰۶۵۱-
S-1559	this is on the fifth of January.
T-1559	میشود پنج ###م ژانویه .
H-1559	این پنج ###م ژانویه است . -۰.۷۹۰۴۷۲۶۸۶۲۹۰۷۴۱-
D-1559	این پنج ###م ژانویه است . -۰.۷۹۰۴۷۲۶۸۶۲۹۰۷۴۱-

استفاده از Tensorboard برای بررسی فرایند آموزش

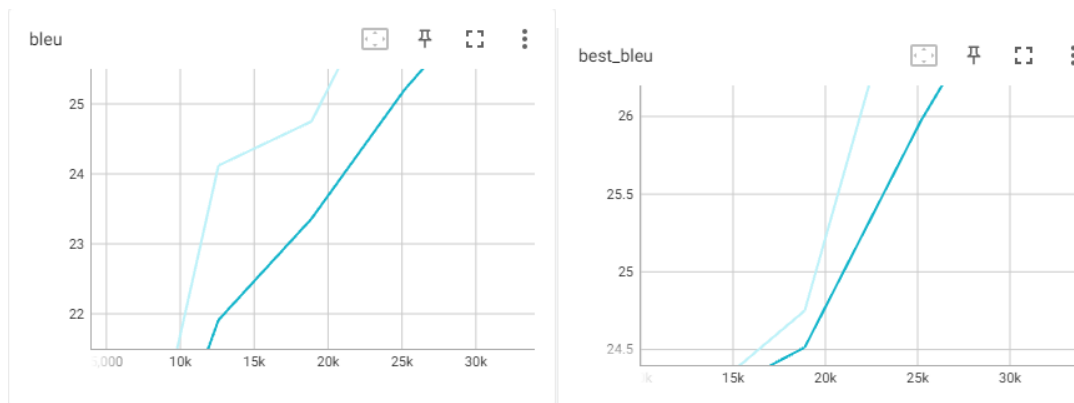
Freeze شده

Loss به شکل زیر می باشد:



شکل ۵ loss مدل freeze شده و استفاده از Bert embedding

نمودار BLEU برای دادگان ارزیابی به شکل های زیر می باشد:



شکل ۶ Bleu مدل freeze شده و استفاده از Bert embedding

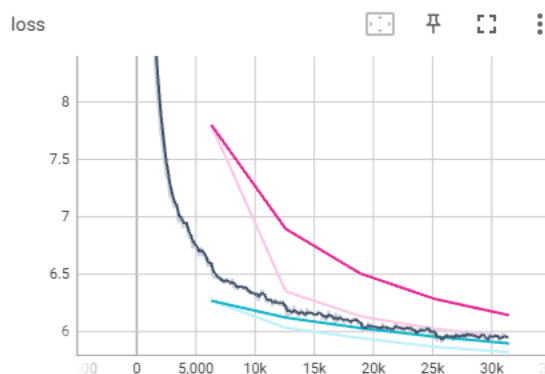
همچنین مقدار نهایی Bleu برای داده‌های ارزیابی به صورت زیر می‌باشد:

Generate test with beam=5: BLEU4 = 17.46, 42.8/22.6/13.1/7.3 (BP=1.000, ratio=1.140, syslen=156777, reflen=137543)

که مدل آنچنان خوب نیست.

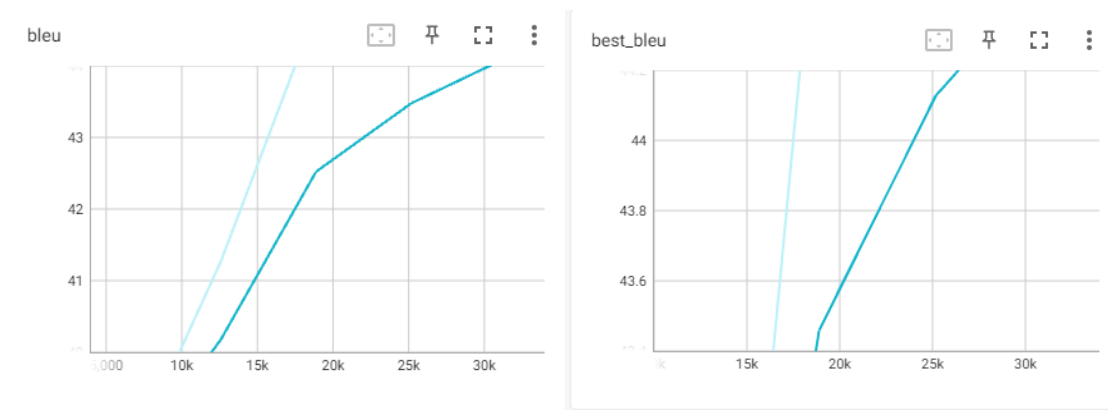
Freeze نشده

Loss به شکل زیر می‌باشد:



شکل ۷ loss مدل freeze نشده و استفاده از Bert embedding

نمودار BLEU برای دادگان ارزیابی به شکل‌های زیر می‌باشد:



شکل ۸ مدل Bleu freeze نشده و استفاده از Bert embedding

مقدار نهایی Bleu برای داده‌های ارزیابی به صورت زیر می‌باشد:

Generate test with beam=5: BLEU4 = 35.53, 62.9/41.9/29.5/20.5 (BP=1.000, ratio=1.012, syslen=139244, reflen=137543)

که با توجه به این نمره مشخص است که به شدت از حالت قبلی مدل بهتر و ترجمه‌های قابل قبولتری دارد و با توجه به فریز نشدن مدل این انتظار می‌رفت.