

۱. برای ساخت شبکه پیش‌بینی کننده خود از یکی از کتابخانه‌های Torch یا TensorFlow استفاده کنید. از شبکه‌ی LSTM و CNN استفاده کرده و داده‌ها را به عنوان ورودی به آن‌ها بدهید و سپس موارد خواسته شده را همراه با توضیح مختصر در سند خود ذکر کنید.
۲. پیش‌پردازی بر روی داده انجام دهید (در صورت لزوم) و دلیل هر پیش‌پردازش را ذکر کنید.

- حذف کلمات توقف

به‌منظور حذف کلمات توقف^۳ و علائم نگارشی^۴، از فایل کلمات توقف فارسی^۵ استفاده شده است اما لازم است در ابتدا پردازش مختصری صورت پذیرد و برخی کلمات مانند کلمات «دیر» یا «زود» از این لیست حذف شوند. زیرا در تعیین میزان رضایت کاربر نقش دارند. این کاراکترها و کلمات، به‌دلیل تمیز کردن مجموعه داده و آماده‌سازی آن برای ساخت مدل، حذف شده‌اند.

```
path = "/content/drive/MyDrive/persian_stopwords.txt"
stop_words = []
f = open(path, "r", encoding='utf-8-sig')
for x in f:
    stop_words.append(x.rstrip("\n"))

train_data['comment'] = train_data['comment'].apply(lambda x: ' '.join([word for word in x.split() if word not in (stop_words)]))
validation_data['comment'] = validation_data['comment'].apply(lambda x: ' '.join([word for word in x.split() if word not in (stop_words)]))
```

- تبدیل متن انگلیسی به فارسی

برخی از نظرات ثبت‌شده در این مجموعه داده یا بخشی از آن‌ها به زبان انگلیسی نوشته شده است. به‌منظور تبدیل این جملات به متن فارسی، از یک تابع استفاده شده است که در ابتدا بخشی از نظرات که با الفبای غیرفارسی نوشته شده را جدا نموده و سپس با استفاده از ماژول تشخیص زبان کتابخانه‌ی googletrans، اقدام به تشخیص زبان متن می‌کند. در صورتی‌که زبان متن مذکور انگلیسی باشد، با استفاده از تابع translate متن ترجمه شده و بازگردانده می‌شود. این کار برای پردازش تمامی نظرات درون مجموعه داده و استفاده کردن از آن‌ها در ساخت مدل انجام شده است.

³ Stop Words

⁴ Punctuation

⁵ <https://github.com/kharazi/persian-stopwords>

```
import re
from googletrans import Translator

def english_handle(text):
    english_part = ""
    farsi_part = ""

    translator=Translator()

    for word in text.split():
        z = re.match("[A-Za-z]+", word)
        if z:
            english_part += word + ' '
        else:
            farsi_part += word + ' '
```

```
if english_part != "":
    a = translator.detect(english_part)
    if a.lang == 'en':
        text=(english_part)
        destination_language = {
            "Persian": "FA",
        }

        for key, value in destination_language.items():
            english_part = translator.translate(text, dest=value).text

new_text = farsi_part + english_part

return new_text
```

• نرمال سازی متن

به منظور نرمال سازی متن از کتابخانه‌ی پارسی‌وار^۶ استفاده شده است. در صورتی که مقدار پارامتر `pinglish_conversion_needed` نرمال ساز برابر `true` قرار گیرد، متن‌های فارسی که با الفبای انگلیسی نوشته شده باشند به شکل املاي فارسی بازگردانده می‌شوند. هم‌چنین از `stemmer` برای ریشه‌یابی کلمات استفاده شده است. در پایان اگر کلمه‌ای حاوی کاراکتر غیرفارسی باشد، حذف می‌شود. در برخی از کلمات نیز بعضی از حروف چندین بار تکرار شده‌اند. به منظور یکسان سازی این کلمات، تکرارهای اضافی حروف حذف شده‌اند. به دلیل ایجاد برخی نیم‌فاصله‌های نادرست بین کلمات، نیم‌فاصله‌های ایجادشده توسط نرمال ساز با فاصله جایگزین شده‌اند.

```
from parsivar import Normalizer
from parsivar import Tokenizer
from parsivar import FindStems

my_normalizer = Normalizer(pinglish_conversion_needed=True)
my_tokenizer = Tokenizer()
my_stemmer = FindStems()

# second_normalizer = Normalizer(statistical_space_correction=True)

def preprocess(text):
    tokens = my_tokenizer.tokenize_words(my_normalizer.normalize(text))
    stemmed_tokens = []
    for token in tokens:
        stemmed_tokens.append(my_stemmer.convert_to_stem(token))

    clean_text = "".join([word + " " for word in stemmed_tokens])
    clean_text = re.sub('[a-zA-Z]+', ' ', clean_text)
    clean_text = clean_text.replace('\u200c', ' ')

    return clean_text
```

^۶ Parsivar

```
def remove_consec_duplicates(s):
    new_s = ""
    prev = ""
    for c in s:
        if len(new_s) == 0:
            new_s += c
            prev = c
        if c == prev:
            continue
        else:
            new_s += c
            prev = c
    return new_s

train_data['comment']=train_data['comment'].apply(preprocess)
validation_data['comment']=validation_data['comment'].apply(preprocess)

train_data['comment']=train_data['comment'].apply(remove_consec_duplicates)
validation_data['comment']=validation_data['comment'].apply(remove_consec_duplicates)
```

از تعبیه کلمه‌های^۷ word2vec و fastText^۸ برای وزن‌های اولیه‌ی شبکه‌ی خود استفاده کنید. (شما برای هر کلمه‌ی نیاز دارید تا برداری از اعداد را جایگزین آن کنید، تعبیه کلمه برای همین عملیات استفاده می‌شود)

با استفاده از فایل 61.zip می‌توان ماتریس تعبیه کلمات از پیش آموزش یافته را به‌دست آورد.

```
import os
import zipfile
with zipfile.ZipFile('/content/drive/MyDrive/61.zip', 'r') as zip_ref:
    zip_ref.extractall('/tmp/pre_embedding')

embeddings_index = {}
f = open('/tmp/pre_embedding/model.txt', encoding = "ISO-8859-1")
for line in f:
    word, coefs = line.split(maxsplit=1)
    coefs = np.fromstring(coefs, "f", sep=" ")
    embeddings_index[word] = coefs

f.close()

print('Found %s word vectors.' % len(embeddings_index))

max_length = 100
word_index = tokenizer.word_index

embedding_matrix = np.zeros((len(word_index) + 1, max_length))
for word, i in word_index.items():
    embedding_vector = embeddings_index.get(word)
    if embedding_vector is not None:
        # words not found in embedding index will be all-zeros.
        embedding_matrix[i] = embedding_vector
```

سپس لایه embedding ایجاد می‌شود:

^۷ Word embedding

^۸ <http://vectors.nlpl.eu/repository/>

```
embedding_layer = Embedding(input_dim=len(word_index) + 1,
                             output_dim=max_length,
                             weights=[embedding_matrix],
                             input_length=max_length,
                             trainable=False)
```

۳. مدل خود را آموزش داده و معیار F1^۹ را گزارش دهید. معیارهایی از جمله صحت^{۱۰}، دقت^{۱۱} و فراخوانی^{۱۲} را برای مدل خود ذکر کنید.

شما باید مدل خود را برای چند دوره^{۱۳} آموزش دهید. در هر دوره مدل شما تمامی مجموعه‌ی داده را مشاهده می‌کند و در جهت بهتر شدن مدل (بهینه‌سازی) تلاش می‌کند. (برای مثال ۱۰ دوره)

مقدار معیارهای نام‌برده برای ارزیابی مدل CNN و LSTM با استفاده از مجموعه “dev” در جدول زیر آورده شده است.

مدل معیار	F1	صحت	دقت	بازخوانی	هزینه
CNN	۰.۸۲	۰.۸۱۵	۰.۸۲	۰.۸۲	۰.۵۵۷
LSTM	۰.۷۸	۰.۸۱۵	۰.۷۸	۰.۷۸	۰.۵۵۷

با توجه به مقادیر جدول فوق، مدل CNN نسبت به LSTM اندکی بهتر است.

197/197 [=====] - 1s 4ms/step - loss: 0.5572 - accuracy: 0.8156
Test accuracy : 0.8156265020370483

	precision	recall	f1-score	support
0	0.81	0.82	0.82	3148
1	0.82	0.81	0.82	3149
accuracy			0.82	6297
macro avg	0.82	0.82	0.82	6297
weighted avg	0.82	0.82	0.82	6297

^۹ Macro-Averaged F-Score

^{۱۰} Accuracy

^{۱۱} Precision

^{۱۲} Recall

^{۱۳} Epoch

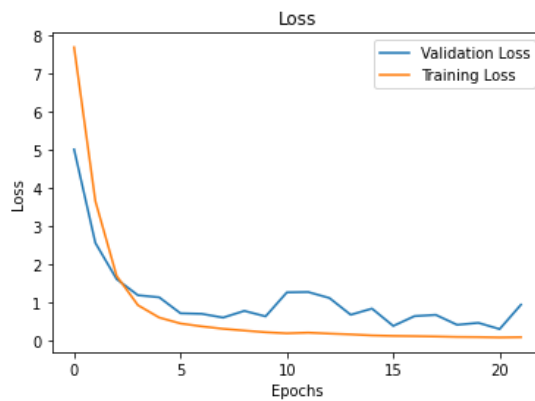
```
197/197 [=====] - 1s 3ms/step - loss: 0.9134 - Accuracy: 0.7818
Test accuracy : 0.7818008661270142
```

	precision	recall	f1-score	support
0	0.78	0.78	0.78	3148
1	0.78	0.78	0.78	3149
accuracy			0.78	6297
macro avg	0.78	0.78	0.78	6297
weighted avg	0.78	0.78	0.78	6297

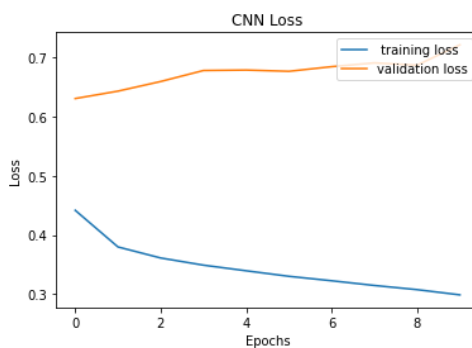
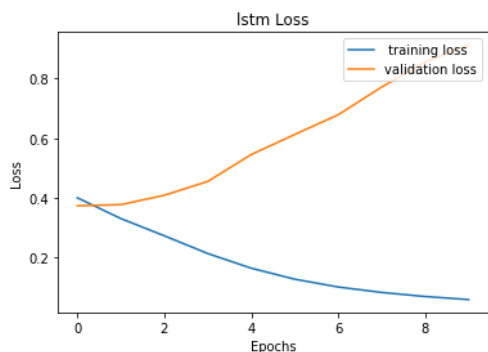
۴. نمودارهای training loss و validation loss را ترسیم کنید و آن را تحلیل کنید.

پس از پایان هر دوره، مقدار loss شبکه‌ی خود در زمان train و در زمان validation را ذخیره کنید. پس از آموزش مدل برای مثلاً ۱۰ دوره، مقادیری دوره‌های مختلف را بصورت نمودار ترسیم نمایید. مثالی از این نمودار آورده شده است که مدل پس از تقریباً ۳ دوره به اصطلاح overfit شده است. (محل تقاطع training loss و validation loss).

با توجه به موارد گفته شده، بهترین مدل خود را ذخیره کرده و برای تست مدل از آن استفاده کنید.

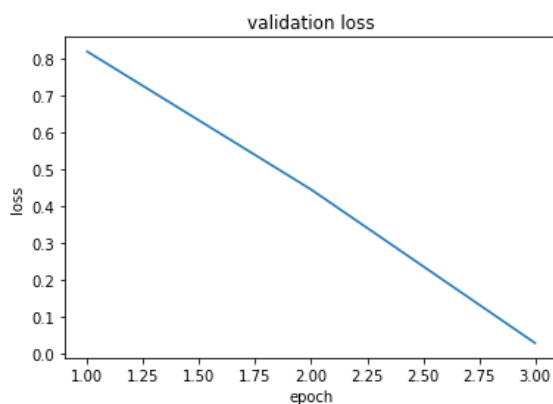


با توجه به نمودارهای training loss و validation loss و عدم تقاطع آن‌ها با یکدیگر، مدل از ابتدا overfit شده است. این حالت برای هر دو مدل رخ داده است.



۵. از مدل پیش آموزش دیده شده، مانند mBERT و XLM-Roberta به عنوان مدل خود استفاده و آن را آموزش دهید. همانند بخش قبل معیار F1 را گزارش و نمودارهای validation loss و training loss را ترسیم کنید. مقایسه‌ای بین مدل‌های قبل و مدل جدید انجام دهید. به‌منظور استفاده از این دو مدل، از کتابخانه‌ی simpletransformers استفاده شده است. نتایج مدل mbert به‌صورت زیر است:

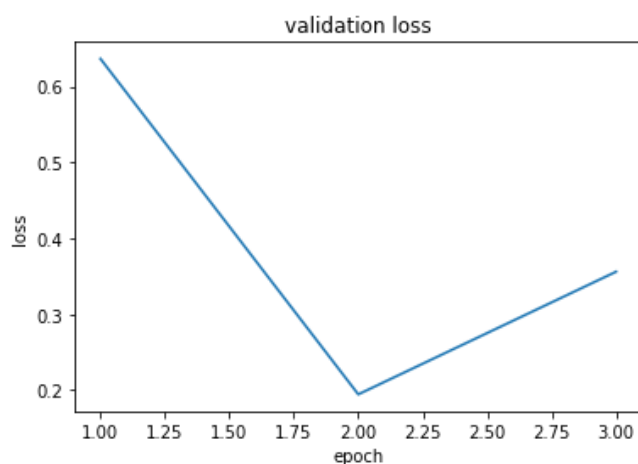
F1	Val_loss	accuracy
0.86	0.372	0.86



```
Epoch 3 of 3: 100% ██████████ 3/3 [50:44<00:00, 1014.25s/it]
Epochs 0/3. Running Loss: 0.8181: 100% ██████████ 7088/7088 [16:45<00:00, 7.31it/s]
Epochs 1/3. Running Loss: 0.4452: 100% ██████████ 7088/7088 [16:50<00:00, 7.30it/s]
Epochs 2/3. Running Loss: 0.0302: 100% ██████████ 7088/7088 [16:42<00:00, 7.57it/s]
INFO:simpletransformers.classification_model: Training of bert model complete. Saved to outputs/.
(21264, 0.4118813318673429)
```

همچنین نتایج مدل xlm_roberta به صورت زیر می باشد:

F1	Val_loss	accuracy
0.863	0.35	0.863



با توجه به اعداد به دست آمده، مدل های از پیش آموزش یافته، نتیجه بهتری دارند.

نکات تکمیلی

- برای درگیر نشدن با کانفیگ ها پیشنهاد می شود از google colab استفاده کنید.
- برای سریع شدن محاسبات پردازش را بر روی GPU و بر روی colab انجام دهید.
- لطفا سند پروژه را حتما در سامانه ی [کوئرا](#) ارسال کنید.
- لطفا پاسخ های خود را در سند پروژه نوشته و در قالب یک فایل PDF ارسال کنید.
- نام سند ارسالی {-Name Family}-Final-{-Name Family} (یک نفر از اعضای گروه ارسال کند)
- برای تحویل پروژه لازم است کدی نوشته شود که یک دیتاست را دریافت و دقت و F1 را بر روی بهترین مدل ذخیره شده توسط شما گزارش می کند. مجموعه داده تست در روز تحویل در اختیار شما قرار خواهد گرفت.
- تمامی فایل های مورد نیاز این تمرین در این [لینک](#) قابل دسترس است.