



## موضوع : تشخیص پیام های تبلیغاتی

ساناز الهیاری و مسعود جانفشان

استاد راهنما : دکتر میلاد سلطانی

دانشکده فنی و مهندسی، گروه کامپیوتر

هدف طراحی برنامه ایست که پیام های تبلیغاتی را از پیام های غیر تبلیغاتی تشخیص دهد

برای این برنامه از 4 مدل میتوان استفاده کرد

با تست های گرفته شده از 4 مدل گزارش زیر بدست آمده :

	accuracy	precision	recall	f1-score
MultinomialNB Model	0.962332	1.000000	0.720000	0.837209
Custom-Vec-Embedding Model	0.981166	0.970803	0.886667	0.926829
Bidirectional-LSTM Model	0.975785	0.936170	0.880000	0.907216
USE-Transfer learning Model	0.981166	0.944828	0.913333	0.928814

در چهار مدل بالا همگی از ساختار شبکه ی پرسپترون پیروی میکنند و تفاوت در ابعاد لایه ها و تعداد نرون های استفاده شده (معماری شبکه ) در این مدل هاست

در جدول بالا از چهار مولفه اصلی برای مقایسه و ارزیابی مدل ها استفاده شده :

accuracy : میزان دقت هر مدل در تشخیص درست نمونه ها به نسبت کل داده ها ( مثبت – مثبت و منفی – منفی )

Precision : درصد کل نمونه های واقعی مثبت / تعداد نمونه های تشخیص داده مثبت

Recall : برعکس Precision ارزیابی میکند ( درصد تعداد نمونه های تشخیص داده مثبت / درصد کل نمونه های واقعی مثبت )

F1-score : در نهایت این مدل ارزیابی هم precision و هم recall را در نظر میگیرد و فرمول آن به شکل زیر است ( جامع ترین معیار برای ارزیابی ) 
$$F1\text{-score} = 2 \times ( \text{Precision} * \text{Recall} ) / ( \text{Precision} + \text{Recall} )$$

\*در نهایت با محاسبه این چهار معیار برای هر مدل گزارشی مانند جدول بالا تهیه شده است \*

## الگوریتم مدل Custom-vec-embedding

این مدل از لایه های مختلف با ساختار پرسپترونیک تشکیل شده :

- لایه TextVectorization برای تبدیل متن خام به یک فرمت عددی مناسب برای ورود به شبکه عصبی استفاده می شود. این لایه متن را توکنیزه کرده، آن را استاندارد می کند (با تبدیل به حروف کوچک و حذف علائم نگارشی) و دنباله هایی از اعداد صحیح که کلمات را نمایندگی می کنند را خروجی می دهد.

- لایه Embedding مسئول یادگیری نمایندگی های معنایی برای کلمات در دنباله های ورودی است. این لایه دنباله های صحیح شده از TextVectorization را می گیرد و هر کلمه را به یک فضای برداری چگال با اندازه ثابت (در اینجا ۱۲۸) نگاشت می کند.

\* دو لایه بالا بر اساس تعداد کلمات یکتا در دیتا ست ورودی و متوسط تعداد کلمات در هر پیام تنظیم میشوند

- معماری مدل شامل یک لایه ورودی، لایه TextVectorization، لایه Embedding، لایه GlobalAveragePooling1D، لایه Flatten، یک لایه چگال با فعال سازی ReLU و یک لایه چگال نهایی با فعال سازی سیگموئید برای دسته بندی دودویی (اسپم یا نه اسپم) است.

توضیح GlobalAveragePooling1D :

لایه GlobalAveragePooling1D یک لایه پردازش میانگین سراسری است که بر روی یک محور اعمال می شود. در اینجا، ما از GlobalAveragePooling1D بر روی محور طول دنباله های ورودی ( $axis=1$ ) استفاده می کنیم. این لایه به ازای هر نمونه (سطر) از ورودی، میانگین مقادیر در هر بعد از دنباله را محاسبه می کند.

برای درک بهتر، فرض کنید که دنباله های ورودی شامل تعدادی بردار هستند. برای هر بردار، GlobalAveragePooling1D مقدار میانگین اعداد آن بردار را محاسبه می کند. این مقدار میانگین به جای یک بردار چگال، یک عدد است

مزیت اصلی از استفاده از این لایه این است که باعث کاهش ابعاد داده می‌شود. به جای این که با یک بردار چگال برای هر دنباله کار کنیم، میانگین‌گیری از تمام ابعاد بردارها باعث ایجاد یک بردار با ابعاد کمتر می‌شود. این کاهش ابعاد می‌تواند کمک کننده باشد تا مدل سریع‌تر آموزش ببیند و از بیش‌برازش (overfitting) جلوگیری کند.

- در نهایت لایه آخر قرار دارد که یک نورون تعریف شده و با تابع فعالساز sigmoid تنظیم شده است ( خروجی بازه ی 0 تا 1 می دهد )

از بین مدل ها Custom-Vec-Embedding را انتخاب می کنیم چون بهترین کارایی و دقت را دارد

در کد کتابخانه های زیر استفاده شده :

Numpy

Pandas

Tensorflow & Keras

sklearn

برای شروع کار یک دیتا ست به فرمت csv ( به نام train) داریم که حاوی حدودا پنج هزار نمونه پیام است که برای آموزش مدل تهیه شده است

با استفاده از دستور read\_csv در کتابخانه pandas فایل را میخوانیم که به عنوان آموزش انتخاب کردیم

(train) را میخوانیم و در متغیری به اسم df میریزیم (آدرس فایل را میدهم)

```
df = pd.read_csv("D:/Machine L projects/train.csv",encoding='latin-1')
```

برای پیش پردازش فایل ورودی ، ستون های اضافه را حذف کرده و نام ستون هایی که حاوی اطلاعات مورد نیاز برنامه هستند را به text و label عوض میکنیم (با استفاده از دستور rename)

```
df = df.drop(['Unnamed: 2','Unnamed: 3','Unnamed: 4'],axis=1)
df = df.rename(columns={'v1':'label','v2':'Text'})
```

سپس ستونی ایجاد میکنیم که پیام هایی که تبلیغاتی هستن را یک و غیر تبلیغاتی را صفر قرار میدهد

```
df['label_enc'] = df['label'].map({'ham':0,'spam':1})
```

ستون text حاوی پیام ها و ستون label\_enc که حاوی صفر و یک است را به آرایه numpy تبدیل میکنیم که به فرمت مطلوب برای ورود به مدل ، تبدیل میشه و به ترتیب در متغیر X و y میریزیم

```
X, y = np.asarray(df['Text']), np.asarray(df['label_enc'])
```

با دستور dataframe جدولی درست میکنیم با دو ستون

یک ستون با نام text که حاوی پیام ها است و ستون بعدی با نام label که حاوی صفر و یک است

```
new_df = pd.DataFrame({'Text': X, 'label': y})
```

حال دو متغیر جدید به نامهای X\_train و y\_train تعریف می‌شوند که به ترتیب برای نگهداری داده‌های متنی new\_df و برچسب‌های مربوط به آن‌ها استفاده می‌شوند

```
X_train, y_train = new_df['Text'], new_df['label']
```

کلمات هر پیام را با استفاده از تابع split جدا کرده سپس تعداد کلمات را بدست می‌آورد و با تابع sum مجموع تعداد کلمات همه ی پیام ها را محاسبه میکند سپس تقسیم بر تعداد پیام ها میکند آن را رند میکند تا میانگین کلماتی که در پیام ها وجود دارد را بدست آورد و در متغیر avg\_word\_len میریزد

```
avg_words_len = round(sum([len(i.split()) for i in df['Text']])/len(df['Text']))
```

با استفاده از تابع set لیستی از کلمات یکتا درست می‌کنیم و تعداد کلمات را در متغیر total\_words\_length می‌ریزیم

```
s = set()
for sent in df['Text']:
    for word in sent.split():
        s.add(word)
total_words_length=len(s)
```

در اینجا تابع کمکی به اسم `fit_model` برای آموزش مدل طراحی شده که پارامترهای ورودی آن ورودی و خروجی (برچسب های آموزش) یعنی `X_train` و `y_train` هستند و در پارامتر `epochs` تعداد تکرار آموزش را مشخص میکنیم

```
def fit_model(model, epochs, X_train=X_train, y_train=y_train,):
    '''
    fit the model with given epochs, train
    and test data
    '''
    history = model.fit(X_train,
                        y_train,
                        epochs=epochs,
                        )
    return history
```

متغیر `MAXTOKENS` تعیین کننده حداکثر تعداد کلمات مورد پذیرش است که مقدار کلمات یکتا را به آن میدهیم و `OUTPUTLEN` مقدار طول خروجی پیش بینی شده مفید است که میانگین کلمات را در آن مقدار دهی می کنیم.

```
MAXTOKENS=total_words_length
OUTPUTLEN=avg_words_len
```

## پیش پردازش داده های متنی :

از کلاس TextVectorization در TensorFlow کتابخانه Keras برای ایجاد یک بردار از داده های متنی استفاده می کنیم که پیام ها به فرمتی تغییر کنند که الگوریتم های یادگیری عمیق بتوانند با آنها کار کنند (تبدیل به داده های عددی) که پارامترهای آن به صورت زیر است :

از max\_tokens برای تعداد بیشینه توکن ها (کلمات) که باید برای واژه نامه بردار سازی در نظر گرفته شود استفاده میشود

با استفاده از پارامتر standardize='lower\_and\_strip\_punctuation' همه متن ها پیش از بردار سازی یکدست می شوند (تمام حروف به حروف کوچک تبدیل شده و علائم نگارشی حذف می شوند).

پارامتر output\_mode نحوه خروجی دادن داده های بردار سازی شده را مشخص می کند. که در اینجا int به این معناست که خروجی برای هر کلمه عدد صحیح منحصر به فردی خواهد بود

پارامتر output\_sequence\_length طول ثابت خروجی هر بردار را تنظیم می کند. اگر متن طولانی تر از OUTPUTLEN باشد، برش خواهد خورد و اگر کوتاه تر باشد، با صفر پر خواهد شد تا به این طول برسد

در نهایت text\_vec.adapt(X\_train) روی مجموعه داده آموزشی X\_train عمل می کند تا لایه بردار سازی را برای تبدیل دقیق و کارآمد داده های متنی آینده آماده سازد. بعد از اجرای این متد، لایه TextVectorization می تواند متون جدید را با استفاده از همین واژه نامه بردار سازی کند، که به آن کمک می کند هر کلمه ای را به معادل عددی اش (که یک اندیس است) تبدیل کند.

```
text_vec = TextVectorization(
    max_tokens=MAXTOKENS,
    standardize='lower_and_strip_punctuation',
    output_mode='int',
    output_sequence_length=OUTPUTLEN
)
text_vec.adapt(X_train)
```



لایه‌های Embedding معمولاً به کار می‌روند تا داده‌های متنی که به صورت اندیس‌های عددی هستند، را به بردارهای چند بعدی تبدیل کنند. این می‌تواند به مدل کمک کند تا مفاهیم پیچیده و رابطه‌های معنایی بین کلمات را بهتر یاد بگیرد

با استفاده از این لایه ارزش ویژگی‌ها را بدست می‌آوریم که ببینیم کدام ویژگی‌ها برای تشخیص تبلیغاتی بودن یا نبود مناسب است

اول یک نمونه جدید از کلاس Embedding را ایجاد می‌کنیم و نتیجه را در یک متغیر به نام embedding\_layer قرار می‌دهیم

پارامتر input\_dim مشخص می‌کند که تعداد کل کلماتی که لایه Embedding باید پشتیبانی کند چقدر است. در این مورد، به تعداد مشخص شده توسط متغیر MAXTOKENS (کلمات یکتا) اشاره دارد

output\_dim تعداد بعد برداری است که هر توکن به آن تبدیل می‌شود. هر کلمه در نمایش برداری به یک بردار 128 بعدی تبدیل خواهد شد.

embeddings\_initializer تعیین می‌کند که وزن‌های اولیه لایه Embedding چگونه مقداردهی شوند. uniform بدین معنی است که وزن‌ها از یک توزیع یکنواخت در یک بازه مشخص، به طور تصادفی انتخاب می‌شوند.

input\_length طول دنباله‌های ورودی را مشخص می‌کند که باید برابر با output\_sequence\_length باشد تا بتوان آن‌ها را پردازش کرد

```
embedding_layer = layers.Embedding(
    input_dim=MAXTOKENS,
    output_dim=128,
    embeddings_initializer='uniform',
    input_length=OUTPUTLEN )
```

حال یک مدل از شبکه عصبی به نام Custom-Vec-Embedding را با استفاده از کتابخانه Keras در TensorFlow ایجاد می‌کنیم

ابتدا یک لایه ورودی تعریف می‌کنیم  $shape=(1,)$  نشان می‌دهد که هر ورودی دقیقاً یک متن است و  $dtype=tf.string$  مشخص می‌کند که نوع داده‌های ورودی، رشته‌های متنی هستند.

`text_vec` که قبلاً تعریف شده بود را روی لایه ورودی اعمال می‌کند تا ورودی از نوع متنی به اندیس‌های عددی تبدیل شود.

سپس `vec_layer` بر روی خروجی لایه بردارسازی اعمال می‌شود. این خروجی حاوی اندیس‌های عددی است که حالا به بردارهای پیچیده‌تر تبدیل می‌شوند.

و با عبور بردارهای بدست آمده از لایه `Global Average Pooling 1D` ابعاد بردارها را کاهش می‌دهیم (میانگین بردارها را در هر اندیس محاسبه می‌کنیم) و در متغیر `x` میریزیم سپس لایه `Flatten` آن را از ابعاد دوبعدی به یک بردار یک بعدی تبدیل می‌کند تا بتواند به عنوان ورودی به لایه‌های بعدی انتقال داده شود متغیر `x` آپدیت میشود و در آخر از یک لایه تمام اتصال (`layers.Dense`) با 32 نورون و تابع فعال‌سازی `ReLU` عبور می‌کند

لایه خروجی، که تمام اتصال است، با یک نورون معرفی می‌شود و از تابع فعال‌سازی سیگموید استفاده می‌کنیم که برای سنجش خروجی‌های دودویی به کار می‌رود خروجی بین صفر و یک است اگر به صفر نزدیک تر باشد صفر در غیر این صورت یک بدست می‌آید.

```
input_layer = layers.Input(shape=(1,), dtype=tf.string)
vec_layer = text_vec(input_layer)
embedding_layer_model = embedding_layer(vec_layer)
x = layers.GlobalAveragePooling1D()(embedding_layer_model)
x = layers.Flatten()(x)
x = layers.Dense(32, activation='relu')(x)
output_layer = layers.Dense(1, activation='sigmoid')(x)
```

در نهایت یک مدل Keras جدید با ورودی و خروجی تعیین شده‌ای ایجاد می‌کنیم

و مدل را کامپایل می‌کنیم این کار محیط مدل را برای آموزش آماده می‌کند.

پارامتر 'optimizer' تعیین می‌کند که چه الگوریتم بهینه‌سازی برای تغییر وزن‌ها و بیاس‌ها در فرایند آموزش به کار برده شود.

loss تابع زیانی است که در طی آموزش سعی می‌شود کمینه شود (میزان عدم تطابق بین برچسب‌های پیش‌بینی شده توسط مدل و برچسب‌های واقعی).

label\_smoothing=0.5 یک تکنیک است که به هدف بهبود روند تعمیم‌پذیری مدل، سخت‌گیری مدل نسبت به برچسب‌های داده آموزشی را کاهش می‌دهد.

در پایان، مدل با استفاده از داده‌های آموزشی X\_train و برچسب‌های y\_train برای تعداد 5 دوره (epoch) آموزش داده می‌شود. نتایج آموزش در متغیر history\_1 ذخیره می‌شوند

```
model_1 = keras.Model(input_layer, output_layer)

model_1.compile(optimizer='adam', loss=keras.losses.BinaryCrossentropy(
    label_smoothing=0.5))

history_1 = model_1.fit(X_train, y_train, epochs=5)
```

تا اینجا به کد آموزش‌های لازم را داده ایم و آماده دریافت ورودی است

از کاربر می‌خواهیم تصمیم بگیرد که تعدادی پیام را در قالب فایل به عنوان ورودی می‌دهد (ds) یا یک پیام را (ex)

اگر کاربر ex را انتخاب کرد وارد تابع `user_exaple()` می‌شویم و اگر ds را انتخاب کرد وارد تابع

`dataset_result()` می‌شویم و اگر غیر از این دو را وارد کرد دوباره سوال را از او بپرسد

```
while True :
    us = input("Do you want the results of your dataset or do you have another
example? [ex/ds]")
    if us=="ex" :
        user_example()
        break

    elif us=="ds" :
        dataset_result()
        break
    else:
        print("Please enter either ex for an example or ds for a dataset.")
```

تابع `user_example()`

در اینجا، Z به عنوان ورودی از کاربر گرفته شده و سپس به یک لیست از متن‌ها تبدیل شده است. سپس از این

لیست برای ساخت دیتافریم استفاده می‌شود و در نهایت مدل بر روی داده‌های تست پیش‌بینی می‌شود.

در نتیجه، پس از اجرای این خط کد، تمام داده‌ها در `Z_test` نگه داشته می‌شود و قابل استفاده برای ارزیابی مدل

یا پیش‌بینی در آینده است.

پیام را وارد مدل ساخته شده می‌کنیم اگر پیام تبلیغاتی باشد خروجی spam در غیر این صورت ham است

```
def user_example():
    z=(input("Mesal khod ra vared konid :"))
    text_list = [z]
    df_user = pd.DataFrame({'Text': text_list})
    Z_test = df_user['Text']
    p = model_1.predict(Z_test)
    if np.any(np.round(p)==1):
        print('spam')
    else:
        print ('ham')
```

## تابع dataset\_result()

از کاربر می خواهیم آدرس فایل مورد نیاز را وارد کند و با استفاده از دستور read آن را میخوانیم

برای پیش پردازش فایل ورودی ، ستون های اضافه را حذف کرده و نام ستونی را که حاوی اطلاعات مورد نیاز

برنامه است را به text عوض میکنیم (با استفاده از دستور rename)

ستون text حاوی پیام ها است را به آرایه numpy تبدیل میکنیم که به فرمت مطلوب برای ورود به مدل ، تبدیل

میشود و در متغیر U میریزیم سپس با دستور dataframe جدولی درست میکنیم با یک ستون به نام text که

حاوی پیام های مورد نظر است

متغیر جدیدی به نام X\_test تعریف می شوند که برای نگهداری داده های متنی new\_df2 استفاده میشود

پیام ها را وارد مدل ساخته شده می کنیم اگر پیام تبلیغاتی باشد خروجی یک در غیر این صورت صفر است

```
def dataset_result():
    df_2 = pd.read_csv(input("Adress file ra vared konid :"),encoding='latin-1')
    #D:/Machine L projects/test.csv

    df_2 = df_2.drop(['Unnamed: 2','Unnamed: 3','Unnamed: 4'],axis=1)
    df_2 = df_2.rename(columns={'v2':'Text'})

    U = np.asanyarray(df_2['Text'])
    new_df_2 = pd.DataFrame({'Text': U})
    X_test = new_df_2['Text']

    predictions = model_1.predict(X_test)

    result_df = pd.DataFrame({
        'Text': df_2['Text'],
        'Prediction': np.round(predictions).flatten()
    })
    result_df['Prediction'] = result_df['Prediction'].replace({0: 'Not Spam ', 1:
'Spam'})

    result_df.to_csv('output.csv', index=False)
```

در نهایت در ذخیره سورس فایلی به فرمت CSV ایجاد میشود که شامل دو ستون Text و Prediction (پیش بینی) است و اگر خروجی مدل 0 باشد Not spam و اگر 1 باید Spam به ما نمایش داده میشود (درستون prediction)

## منابع

[geeksforgeeks.org](https://www.geeksforgeeks.org)

[kaggle.com](https://www.kaggle.com)

[chatgpt.openai](https://chatgpt.openai.com)

[stackoverflow.com](https://stackoverflow.com)