

گزارش تمرین تشخیص ایمیل اسپم

نیکتا گوهری صدر 9513433

در ابتدا کتابخانه های مورد نظر ایمپورت می شوند

سپس فایل ها خوانده میشوند که در اینجا به جای کلمه ham از okay استفاده شده است.

سپس داده ها پیش پردازش می شوند که این پیش پردازش شامل جداسازی کلمات یا tokenize، حذف stop word ها یا ایست واژه ها، ریشه یابی کلمات یا stemming و همچنین چک کردن اینکه بعد از ریشه یابی طول کلمه کمتر از ۲ نباشد و حذف کاراکترهای غیر فارسی مانند حروف انگلیسی، علامت گذاری ها یا punctuation، اعداد و غیره میباشد.

به دلیل این که داده ها جداگانه خوانده شده اند پیش پردازش داده های اسپم و غیر اسپم جداگانه صورت می گیرد

همچنین پیش پردازش از دو لیست استفاده می کند که لیست اول کلمات مناسب را در خود ذخیره می کند و سپس بعد از ریشه یابی و تمامی مراحل گفته شده، کلمات دوباره جوین شده و به صورت جمله در لیست نهایی ذخیر میشود.

این پیش پردازش هم روی داده های train و هم test انجام میشود که نمونه ی نتیجه ی آن در عکس زیر مشخص است.

عکس اول قبل از پیش پردازش است:

0	
0	با سلام و احترام بنده بواسطه پیرو نامه شماره ت...
1	«لا یحب الله جهر بالسوء من القول الا من ظلم»...
2	سلام. در مورد این که آموزش و پرورش مدرک آرها ...
3	سلام چه خوبه که شما احساس مسئولیت دارید و خوا...
4	سلام من مترجمی زبان انگلیسی خندم البته نه سنا...

عکس دوم بعد از پیش پردازش است:

0	سلام احترام بدینوسیله پیرو نامه شماره تاریخ اه...
1	بخت الله جهر بالسوء القول الا ظلم جناب آقا حرم ...
2	سلام آموزش پرورش مدرک آرفا پرسنل رسمی گنه نوی ...
3	سلام خویه احساس مسئولیت جواب دین خوام رانندگی...
4	سلام منرجمی زبان انگلیسی خوندن سال دلائل دست ت...

پس از پیش پردازش معیار خی دو را پیاده سازی می کنیم که برای این کار باید داده های اسپم و ok را ابتدا به هم بچسبانیم و لیست لیبل را هم درست کنیم.

سپس با استفاده از تابع `countvectorizer` داده ها را وکتورایز می کنیم و بهترین ویژگی ها را انتخاب میکنیم که برابر است با k تا ویژگی ای که خی دو بیشترین امتیاز را میدهد. که در این جا $k=500$ قرار داده شده است. که نتیجه ی آن را در عکس زیر مشاهده میکنید.

```
In [15]: 1 train_vector.shape
Out[15]: (500, 13416)

In [16]: 1 type(train_vector)
Out[16]: numpy.ndarray

In [17]: 1 test_vector.shape
Out[17]: (400, 13416)

In [18]: 1 type(test_vector)
Out[18]: numpy.ndarray

In [19]: 1 train_vector_chi2.shape
Out[19]: (500, 500)

In [20]: 1 test_vector_chi2.shape
Out[20]: (400, 500)
```

برای استفاده از الگوریتم knn نیاز به یک معیار شباهت داریم که ابتدا معیار cosine similarity و سپس knn را پیاده سازی می کنیم.

نتیجه ی این الگوریتم بدون معیار خبی دو در عکس زیر مشخص است.

[illegible]

حال می‌خواهیم همین الگوریتم را با معیار خ‌ی دو بررسی کنیم که نتیجه در عکس زیر معلوم است.

[illegible]

در مرحله ی بعد از معیار **tfidf** برای **knn** بدون خی دو استفاده میکنیم که اجرای این الگوریتم بسیار زمان بر است اما در نهایت خروجی داده شده و هم در عکس زیر هم در کد مشخص است.

همچنین برای اجرای این الگوریتم به جای داده های به صورت جمله از داده های tokenize شده استفاده شده است.

```
In [31]: 1 y_pred_tfidf = my_knn(3,tokenized_x_concat_train,tokenized_x_concat_test)

In [32]: 1 y_pred_tfidf

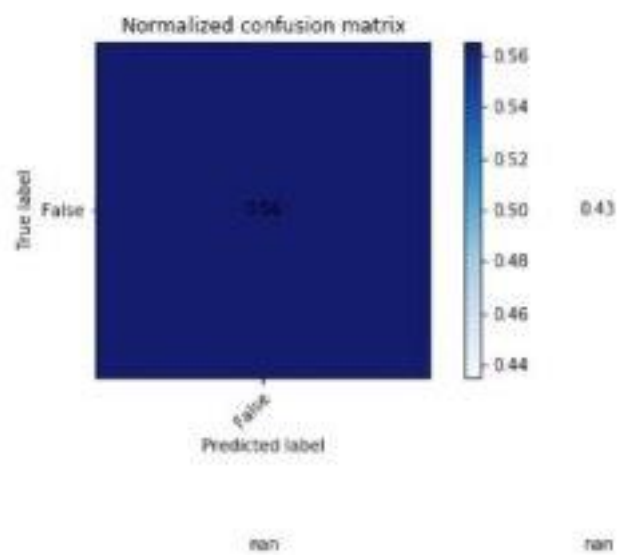
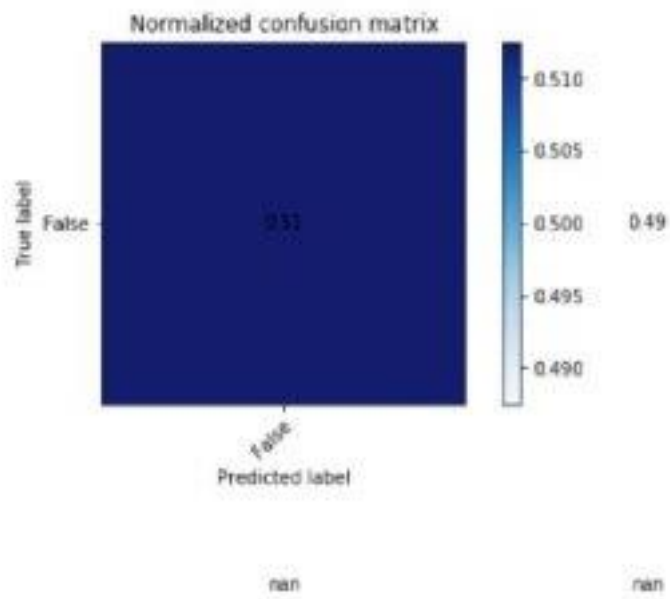
Out[32]: [False,
          False,
          False,
          False,
          False,
          False,
          False]
```

در این مرحله میخواهیم با استفاده از معیار های مختلف مدل هایمان را ارزیابی کنیم.

```
***result of knn with cosine similarity without chi2***
Accuracy: 0.5125
precision: 1.0
recall macro: 0.25625
f1_macro: 0.33884297520661155
recall micro: 0.25625
f1_micro: 0.5125
***result of knn with cosine similarity with chi2***
Accuracy: 0.565
precision: 1.0
recall macro: 0.2825
f1_macro: 0.36102236421725237
recall micro: 0.2825
f1_micro: 0.565
***result of knn with tfidf similarity without chi2***
Accuracy: 1.0
precision: 1.0
recall macro: 1.0
f1_macro: 1.0
recall micro: 1.0
f1_micro: 1.0
/home/nikta/.local/lib/python3.6/site-packages/sklearn/metrics/classification/_classification.py:131: UndefinedMetricWarning: Precision is ill-defined: Labels in y_true not among y_predict labels.
  warnings.warn('Precision is ill-defined: Labels in y_true not among y_predict labels.', UndefinedMetricWarning)
```

همانطور که در عکس مشخص است knn با tfidf بهترین دقت را داشته است و بعد از آن knn با cosine similarity و با feature selection خیی دو نتیجه ی بهتری تا بدون خیی دو داشته است.

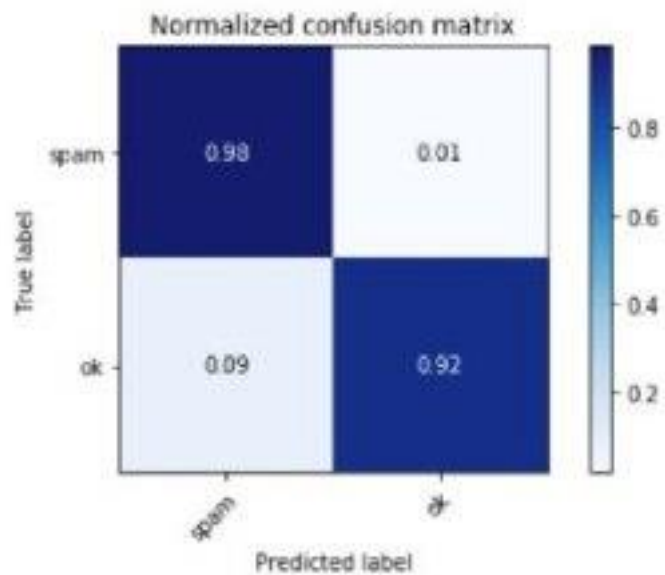
Confusion matrix مدل های نیز در عکس های زیر مشخص است



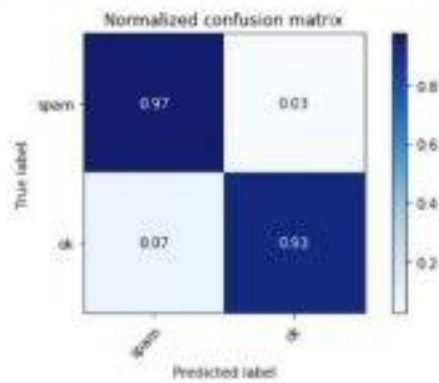

```
In [59]: 1 Y_pred_nb_chi2
```

[illegible]

معیار های ارزیابی و confusion matrix این دو مدل نیز در عکس زیر آمده است.



```
Out[63]: <matplotlib.axes._subplots.AxesSubplot at 8x7f97b5ca5e48>
```



همانطور که میبینیم به طور کلی الگوریتم نایو بیز از knn هم بهتر عمل میکند هم سریع تر. اما مقایسه ی نایو بیز با خی دو و بدون خی دو دشوار است زیرا در برخی معیار ها با خی دو بهتر است و در برخی بدون خی دو.