

به نام خدا
محمد حسین رنگرز - ۹۹۳۶۱۳۰۳۱

1.a:

در این قسمت قرار است با متریک های مختلف اثرگذاری را بررسی کنیم.
اگر داشتن سابجکت را مورد بررسی قرار دهیم داریم:

```
Classification accuracy: 33.1%
```

میتوان مشاهده کرد که دقت به شدت بدی به ما داد و بررسی سابجکت نمیتواند یک متریک خوبی برای پیدا کردن اسپم باشد و فایده ای ندارد.
اگر اثرگذاری و میزان رخداد چند کلمه در کنار هم را بررسی کنیم:

```
Found 57344 spam-indicative 2-grams
```

```
Classification accuracy: 70.3%
```

```
Precision (spam): 72.8%
```

```
Recall (spam): 86.0%
```

باز هم میتوان مشاهده کرد که تأثیر خوبی بر روی پیدا کردن اسپم نداشته است که دقت به حد خیلی خوبی افزایش پیدا کند البته میتوان با N بیشتر بررسی کرد در این مثال صرفاً ۲ گرام بررسی شده است.
اگر مشترک های spam و ham را با یک امتیازی spam در نظر بگیریم:

```
Total spam emails in training: 31134
```

```
Total ham emails in training: 17615
```

```
Classification accuracy: 51.1%
```

```
Precision (spam): 100.0%
```

```
Recall (spam): 23.9%
```

```
F1-score (spam): 38.6%
```

که باز هم میتوان مشاهده کرد درصد بدی به ما میدهد که البته میتوان threshold و اسکور را دستکاری کرد و اون مقدار بهینه شده را بدست آورد شاید بتوان به درصد بهتری برسیم

```
Classification accuracy: 70.6%
```

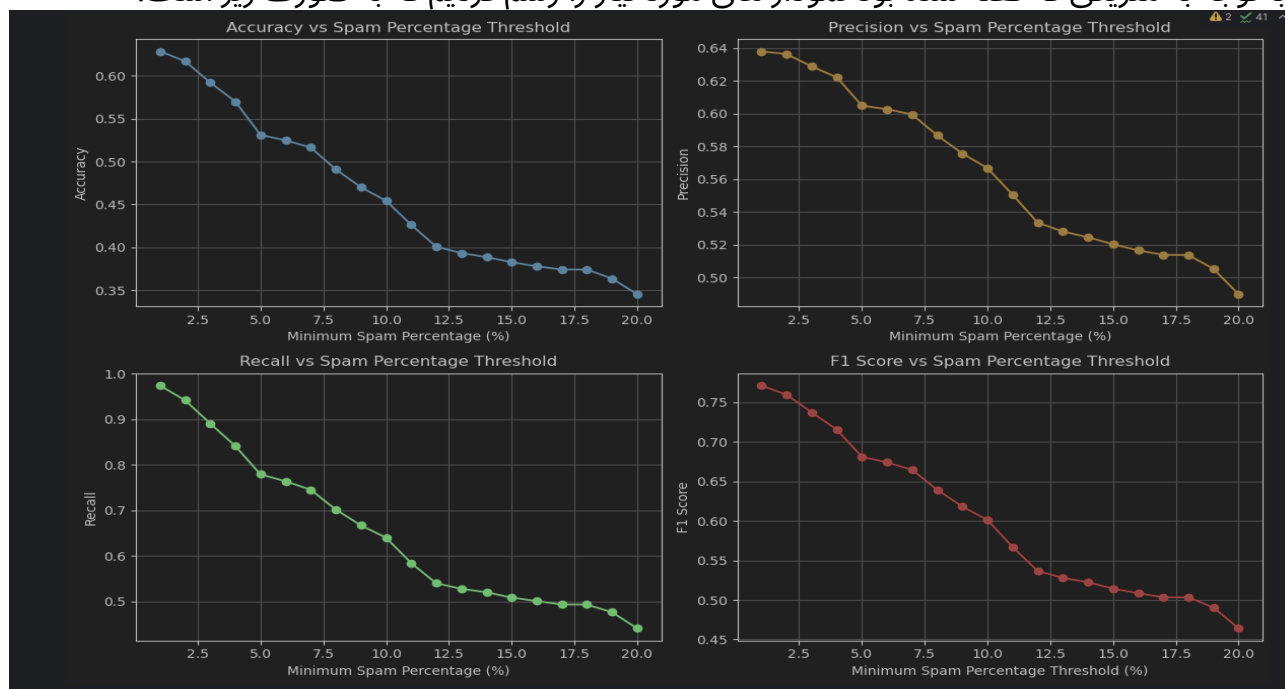
```
Precision (spam): 69.8%
```

```
Recall (spam): 98.5%
```

```
F1-score (spam): 81.7%
```

1.b:

با توجه به متریکی که گفته شده بود نمودار های مورد نیاز را رسم کردیم که به صورت زیر است:



Performance Summary:

Threshold%	Spam Words	Accuracy	Precision	Recall	F1 Score
------------	------------	----------	-----------	--------	----------

1%	1152	0.628	0.638	0.974	0.771
2%	617	0.617	0.636	0.942	0.760
3%	426	0.592	0.629	0.891	0.737
4%	352	0.570	0.622	0.841	0.715
5%	306	0.531	0.605	0.779	0.681
6%	289	0.525	0.603	0.764	0.674
7%	237	0.516	0.600	0.745	0.665
8%	124	0.491	0.587	0.702	0.639
9%	87	0.470	0.576	0.667	0.618
10%	66	0.454	0.567	0.640	0.601
11%	43	0.426	0.551	0.584	0.567
12%	28	0.401	0.533	0.540	0.537
13%	25	0.393	0.528	0.528	0.528
14%	22	0.389	0.525	0.520	0.522
15%	15	0.383	0.520	0.509	0.514
16%	14	0.378	0.517	0.501	0.509
17%	10	0.374	0.514	0.493	0.503
18%	8	0.374	0.514	0.493	0.503
19%	6	0.364	0.505	0.476	0.490
20%	5	0.345	0.490	0.441	0.464

2.c:

همانطور که یکی از متریک هایی که در کلاس مورد بررسی قرار گرفته است بررسی اثر رخداد چند کلمه در کنار هم بود که یکی از روش های این متریک N-gram است. N-gram یک تکنیک پردازش زبان است که ترتیب و ترکیب کلمات را در متن بررسی می کند. در این روش، دنباله های متوالی از N کلمه (مثلاً ۲ کلمه ای Bigram یا ۳ کلمه ای Trigram) استخراج شده و به عنوان ویژگی های معنادار برای تحلیل استفاده می شوند.