



به نام خدا



دانشگاه تهران
دانشکده مهندسی برق و کامپیوتر
هوش مصنوعی

تمرین کامپیوتری سوم

نام و نام خانوادگی	مهسا تاجیک
شماره دانشجویی	810198126
تاریخ ارسال گزارش	29 اردیبهشت

فاز اول : پیش پردازش داده

برای پیش پردازش داده ها از کتابخانه هضم که مخصوص دادگان فارسی است استفاده کردم. در ابتدا تمام punctuation ها را از ستون description حذف کردم که شامل موارد زیر است:

```
Punctuations='\"'#$%&'()*+,-./:;<=>?@[\\]^_`{|}~¡¢£¥¦§¨ª«¬®¯°±²³´µ¶·¸¹º»¼½¾¿'\"'
```

سپس داده ها با تابع normalizer کتابخانه هضم نرمالیزه کردم و فاصله های اضافی را از آن حذف کردم.

سوال 1) stemming : یک تکنیک است که برای استخراج فرم پایه کلمات به حذف نشانه های

جمع ، ضمائر و صفات تفضیلی و عالی را از کلمه میپزدازد که در زبان فارسی شامل موارد زیر است:

["ا"، "ها، ی، ای، ش، ت، م، تر، ترین، ان، ات"]

موتورهای جستجو برای indexing کلمات از stemming استفاده می کنند . به همین دلیل است که یک موتور جستجو به جای ذخیره انواع مختلف کلمه ، می تواند فقط stem ها را ذخیره کند . به این ترتیب ، stemming ساینز شاخص را کاهش می دهد و دقت بازایی را افزایش می دهد.

It is just like cutting down the branches of a tree to its stems. For example, the stem of the words *eating*, *eats*, *eaten* is *eat*.

Lemmatization: این تکنیک مشابه تکنیک قبلی است و خروجی ای که از آن میگیریم lemma

نامیده می شود و در واقع ریشه ی یک کلمه را به ما می دهد.

- هدف هر دو این است که مشتقات یک کلمه را به اصل آن برگرداند تا سرچ کردن راحت تر انجام شود.

فاز دوم : فرآیند مسئله

در این فاز می خواهیم با استفاده از naïve bayes و مدل bag of words ، کتاب هایی که موضوع آن ها نمی دانیم دسته بندی کنیم. برای اینکار باید تعداد هر کلمه را در هر دسته بندی از موضوعات پیدا کنیم. با کنار هم قرار دادن این ها یک ماتریس به ابعاد تعداد کل کلمات(فیچرها) و تعداد کل داکيومنت ها (تعداد سطر های دیتاست) خواهیم داشت. این شمارش را میتوانیم به صورت دستی انجام دهیم یا با استفاده از CountVectorizer در کتابخانه sklearn که کار توکنایز کردن جملات را هم انجام میدهد.

سوال (2)

Prior

در این مسئله احتمالات prior همان احتمال دیده شدن داکيومنت های یک کلاس خاص در کل داکيومنت هاست که در naïve bayes فرض بر این است که احتمال تمام کلاس ها با هم برابر است یعنی تعداد داده هایی که از هر کلاس داریم با هم برابر است و اینجا هم می توانیم با قطعه کد زیر این مسئله را بررسی کنیم. همانطور که میبینیم در کل 2550 داده ی ترین داریم و 6 کلاس داریم که سهم هر کلاس 425 داده است و در واقع احتمال هر کلاس حدودا 0.16 است. برای داده های تست هم این مسئله صدق میکند.

```
1 c1 = c2 = c3 = c4 = c5 = c6 = 0
2 for i in range(len(train_data)):
3     if train_categories[i] == 'مدیریت و کسب و کار':
4         c1+=1
5     elif train_categories[i] == 'رمان':
6         c2+=1
7     elif train_categories[i] == 'کلیات اسلام':
8         c3+=1
9     elif train_categories[i] == 'داستان کودک و نوجوانان':
10        c4+=1
11    elif train_categories[i] == 'جامعه شناسی':
12        c5+=1
13    elif train_categories[i] == 'داستان کوتاه':
14        c6+=1
15
16 pc1 = c1/len(train_data)
17 pc2 = c2/len(train_data)
18 pc3 = c3/len(train_data)
19 pc4 = c4/len(train_data)
20 pc5 = c5/len(train_data)
21 pc6 = c6/len(train_data)
22
23 print(pc1,pc2,pc3,pc4,pc5,pc6)
24 print([len(train_data)])
```

0.16666666666666666 0.16666666666666666 0.16666666666666666 0.16666666666666666 0.16666666666666666 0.16666666666666666
2550

Evidence

احتمال دیده شدن یک کلمه یا فیچر در کل داکيومنتها و تمام کلاس هاست که در این مسئله نیاز به محاسبه ی مستقیم آن نیست.

Likelihood

احتمال دیده شدن یک کلمه در داکيومنت های یک کلاس خاص. در واقع بعد از توکنایز کردن کلمات باید ببینیم هر کدام به تفکیک هر کلاس، در چند تا از داکيومنت ها آمده است.

Posterior

احتمال اینست که یک داکيومنت به شرط دیده شدن فیچرهای تعریف شده ، به یک کلاس تعلق داشته باشد. که در naïve bayes با فرض استقلال ویژگی ها این احتمال برابر با ضرب احتمال های دیده شدن هر کلمه به تفکیک کلاس در احتمال آن کلاس است که مقدار هر کدام بیشتر بود ، داکيومنت به آن کلاس تعلق دارد.

سوال 3) چون سری بر آستانش ز سر صفا نهادی / به صفا و مروه ای دل دگرت چه کار باشد

در مثال بالا معنی کلمه ی صفا در دو جمله متفاوت است که اگر از bigram استفاده کنیم، با در نظر گرفتن صفا و مروه با هم معنی کلمه مشخص میشود و متمایز از صفا در جمله اول خواهد بود.

سوال 4) زمانیکه در بین کلمات به کلمه ی جدیدی برخورد می کنیم که قبلا دیده نشده یا فقط در

یک دسته خاص دیده شود ، در این احتمال شرطی کلاس برای آن کلمه یا ویژگی صفر خواهد شد و در نتیجه احتمال posterior آن صفر خواهد شد و این اتفاق زمانی می افتد که مدل overfit شود یعنی مدل بیش از حد آموزش داده شده و به داده های آموزش چسبیده و قابلیت generalization ندارد.

سوال 5) برای حل مشکل پیش آمده ، از روش additive smoothing میتوان استفاده کرد که در

آن محاسبه احتمال لایکلیهود اندکی متفاوت میشود :

$$P(x_i | C_k) = \frac{(\text{number of points such that } x_i \text{ occurs and class label} = C_k) + \alpha}{\text{number of points where class label} = C_k + \alpha k}$$

Where α can be any real value > 0
 k is the number of class labels.

همانطور که میبینیم در صورت کسر مقدار آلفا جمع می شود و در مخرج $\alpha * k$ که k تعداد کلاس هاست که در مثال ما 6 تاست. و به این شکل از صفر شدن کسر بالا جلوگیری میشود. اگر مقدار الفا را یک قرار دهیم به آن laplace smoothing هم گفته میشود.

سوال (6)

فاز سوم : ارزیابی

سوال (7) غالباً ، یک رابطه معکوس بین precision , recall وجود دارد و مدلی خوب است که هر دوی این مقادیر تا حد تا قبولی بالا باشد که بستگی به مسئله دارد. به طور مثال، جراحی مغز نمونه ای روشن از این tradeoff را ارائه می دهد. یک جراح مغز را در نظر بگیرید که تومور سرطانی را از مغز بیمار خارج می کند. جراح باید تمام سلولهای تومور را از بین ببرد زیرا سلولهای سرطانی باقی مانده باعث بازسازی تومور می شوند. برعکس ، جراح نباید سلولهای سالم مغز را از بین ببرد ، زیرا این کار بیمار را دچار اختلال در عملکرد مغز می کند. جراح ممکن است در ناحیه ای از مغز که خارج می کند ، لیبرال تر باشد تا اطمینان حاصل کند که همه سلول های سرطانی را استخراج کرده است. این تصمیم باعث افزایش recall می شود اما precision را کاهش می دهد. از طرف دیگر ، ممکن است جراح در قسمتی از مغز که برمی دارد محافظه کارتر باشد تا اطمینان حاصل کند که فقط سلول های سرطانی را استخراج می کند. این تصمیم precision را افزایش می دهد اما recall را کاهش می دهد. به عبارت دیگر ، recall بیشتر احتمال از بین بردن سلولهای سالم را افزایش می دهد (نتیجه منفی) و احتمال حذف همه سلولهای سرطانی را افزایش می دهد (نتیجه مثبت). precision بیشتر احتمال از بین بردن سلولهای سالم را کاهش می دهد (نتیجه مثبت) اما احتمال حذف همه سلولهای سرطانی را نیز کاهش می دهد (نتیجه منفی).

معمولا این دو معیار ارزیابی به تنهایی بیان نمی شوند و میتوان threshold های مختلفی در نظر گرفت و در آن ها مقادیر هر دو معیار را محاسبه کرد تا ببینیم برای آن مدل و مسئله کدامیک مناسبتر است.

سوال 8) معیار $F1$ ، میانگین هارمونیک دو معیار $recall$ و $precision$ است که از این جهت اهمیت دارد که بجای اینکه خودمان را درگیر افزایش و تنظیم دو مقدار کنیم فقط سعی میکنیم مقدار این معیار را افزایش دهیم و کار با آن راحت تر است. این معیار ، زمانی که $precision$ و $recall$ مقادیر نزدیک بهم دارند ، تقریبا برابر میانگین آنهاست اما در حالت کلی ، مربع میانگین هندسی آنها تقسیم بر میانگین ریاضی آنهاست.

سوال 9) در این قسمت خواسته شده بود برای هر کلاس معیار $F1$ محاسبه شود که میتوانیم این مقادیر را به یک مقدار کاهش دهیم و تنها یک مقدار گزارش کنیم.

یکی از روش ها برای ادغام $F1$ ها $macro-F1$ است که از میانگین ریاضی بین مقادیر $F1$ برای کلاس های مختلف بدست می آید.

روش های $macro-averaged-recall$ و $macro-averaged-precision$ به ترتیب بین مقادیر تمام $recall$ ها و $precision$ ها میانگین ریاضی می گیرند.

در روش های $weighted-precision$, $weighted-recall$, $weighted-F1$ براساس تعداد سмпل هایی که از هر کلاس داریم مقادیر این معیار ها برای کلاس های مختلف را وزن دهی میکنیم و یک میانگین وزن دار خواهیم داشت.

روش آخر $micro$ است که برابر است با همان $accuracy$:

$$micro-F1 = micro-precision = micro-recall = accuracy$$

سوال 10)

نتایج با استفاده از additive smoothing :

Accuracy برابر است با 82.6٪

```

1 correct = 0
2 for i in range(len(y_pred)):
3     if(y_pred[i] == test_categories[i]):
4         correct+=1
5 accuracy = correct/len(y_pred)
6 print(accuracy)

```

0.8266666666666667

```

Class1 presicion : 0.9253731343283582
Class2 presicion : 0.7530864197530864
Class3 presicion : 0.8732394366197183
Class4 presicion : 0.9375
Class5 presicion : 0.7528089887640449
Class6 presicion : 0.7692307692307693
Class1 Recall : 0.8266666666666667
Class2 Recall : 0.8133333333333334
Class3 Recall : 0.8266666666666667
Class4 Recall : 0.8
Class5 Recall : 0.8933333333333333
Class6 Recall : 0.8

```

```

Class1 F1 : 0.8732394366197183
Class2 F1 : 0.782051282051282
Class3 F1 : 0.8493150684931506
Class4 F1 : 0.8633093525179856
Class5 F1 : 0.8170731707317072
Class6 F1 : 0.7843137254901961

```

macro-f1 : 0.8282170059840066

نتایج بدون استفاده از additive smoothing :

Accuracy برابر است با 78.2٪

```

1 correct = 0
2 for i in range(len(y_pred)):
3     if(y_pred[i] == test_categories[i]):
4         correct+=1
5 accuracy = correct/len(y_pred)
6 print(accuracy)

```

0.7822222222222223

```

Class1 presicion : 0.8955223880597015
Class2 presicion : 0.7073170731707317
Class3 presicion : 0.8208955223880597
Class4 presicion : 0.8676470588235294
Class5 presicion : 0.6956521739130435
Class6 presicion : 0.7567567567567568
Class1 Recall : 0.8
Class2 Recall : 0.7733333333333333
Class3 Recall : 0.7333333333333333
Class4 Recall : 0.7866666666666666
Class5 Recall : 0.8533333333333334
Class6 Recall : 0.7466666666666667

```

```

Class1 F1 : 0.8450704225352113
Class2 F1 : 0.7388535031847134
Class3 F1 : 0.7746478873239437
Class4 F1 : 0.8251748251748251
Class5 F1 : 0.7664670658682634
Class6 F1 : 0.7516778523489932

```

macro-f1 : 0.783648592739325

سوال 11) در استفاده از روش additive smoothing مقدار الفا را برابر با 1 قرار دادم و مقدار accuracy از 78.2% به 82.6% افزایش یافت و مقدار macro-f1 از 78.3% به 82.8% افزایش یافت.

سوال 12)

مشکلی که می تواند وجود داشته باشد اورلپ بین کتگوری هاست مثلاً یک کتاب که در دسته ی داستان کوتاه قرار دارد یا در دسته داستان کودک و نوجوان ، ممکن است این دو، کلمات مشابه زیادی داشته باشند و یک داستان کوتاه کودک و نوجوان داشته باشیم.

مشکل دیگر حذف نکردن کلمات پرتکرار است.

5 کتاب که به اشتباه تشخیص داده شدند که کتگوری های واقعی و پیش بینی شده برای آن ها و توضیحات کتاب آورده شده است.

('رمان', 'داستان کوتاه')

رمان تاریخی «زندانی قلعه قهقهه» روایتگر سرگذشت قهرمانی های شاه «اسماعیل صفوی»، شاه ایران به سالهای ۸۹۲-۹۳۰ ق. است. این داستان از جایی شروع می شود که «مرادبیگ»، نایب الحکومه قلعه قهقهه است. او که مردی تندخو است، زنی زیبا با نام «قمر سلطان» و دختری دلفریب به نام «گل حرم» دارد. در آن دوران مهم ترین زندانی قلعه شاهزاده «اسماعیل میرزا» فرزند شاه «تهماسب دوم» بود که تصمیم داشت با گل حرم ازدواج کند، اما شرایطی پیش آمد که عشق گل حرم را به کینه و دشمنی تبدیل کرد. تا این که اسماعیل میرزا به قزوین می رود و گل حرم از او بی اطلاع می ماند. پس از گذشت سالها، گل حرم نیز به قزوین سفر می کند. در این سفر اتفاقات گوناگونی رخ می دهد که هریک تجربه ای تازه برای گل حرم محسوب می شود. گل حرم تصمیم دارد شاهزاده اسماعیل را بیابد. او شاهزاده را می بیند و علت رفتن او را می یابد. اما پس از آن رخداد های گوناگون دل نگرانی های بسیاری را برای او فراهم می سازد. به راستی سرنوشت گل حرم چه خواهد شد؟

('داستان کوتاه', 'داستان کودک و نوجوانان')

«مطلقاً تقریباً» نوشته لیسا گراف (۱۹۷۴-) نویسنده آمریکایی کتاب های کودکان و نوجوانان است. این کتاب رمانی برای نوجوانان است و از این صحبت می کند که چگونه کشف کنی که هستی و چگونه کاری را که دوست داری انجام دهی. آلبی هرگز زرنگترین شاگرد کلاس نبوده است. هرگز بلندترین نبوده، و هرگز در ورزش هم بهترین نبوده، بهترین هنرمند هم نبوده. در عوض، آلبی یه لیست بلندبالایی از چیزهایی که در آن ها خوب نیست دارد. اما آلبی با پرستار جدیدی آشنا می شود که به او کمک می کند بفهمد در چه چیزهایی خوب است و چطور می تواند به آن کارها مفتخر باشد.

در بخشی از کتاب می خوانیم:

فهمیدم چرا درن و بقیه آن بچه بدجنس ها به آن دختره که آدامس خرسی داشت می گفتند ب - ب - ب - بتسی. چون که گاهی نمی تواند درست و حسابی کلمه ها را بگوید مخصوصاً اول کلمه ها را. مثل ب و ت یا ک. وقتی خانم رز سر کلاس به او گفت که چند خط از روی درس مان بلند بلند بخواند، متوجه این موضوع شدم. پسرها هرهر خندیدند و صورت دختره قرمز شد و صدایش آن قدر آهسته و آرام شد که شنیده نمی شد. تا اینکه بالاخره خانم رز گفت: «متشکرم بتسی، عالی بود».

بتسی زیاد حرف نمی زند.

از کالیستا درباره این موضوع سؤال کردم و کالیستا گفت که شاید بتسی لکنت زبان دارد که حرف زدن را برای آدم سخت می کند.

از بتسی خوشم می آمد. سر ناهار بدون اینکه حتی اسم آدامس بیاورم، به من آدامس خرسی می داد. ما حتی همدیگر را برای کار توی کتابخانه انتخاب کردیم و وقتی خانم رز درباره کارت اطلاعات آنلاین هر کتاب توضیح داد و من گیج شدم، بتسی مسخره ام نکرد. فقط جای درستی را که باید کلیک می کردم، نشانم داد. برایم مهم نبود که بتسی زیاد حرف نمی زد. چون گاهی سخت است که آدم منظورش را بگوید و فکر کردم شاید بیشتر وقتها به هر حال منظورش را می فهمم.

('داستان کودک و نوجوانان', 'کلیات اسلام')

«فاطمه علی است» نوشته علی قهرمانی کتابی درباره زندگی مشترک حضرت زهرا (س) و حضرت علی (ع) است. این اثر می‌تواند الگویی کامل برای زندگی مشترک تمامی زوج های جوان باشد. گذشت ایثار، همدلی و همراهی و عشق و مهرورزی درس‌هایی است که این داستان زندگی به ما می‌دهد:

اول ازدواج‌شان و آغاز راه زندگی مشترک بود. هر دو پیش رسول خدا (ص) آمدند. پیامبر (ص) معلم زندگی بهتر و تدبیر منزل انسان‌ها بود. نوبت که به کارهای خانه رسید، پیامبر (ص) پیشنهاد کرد: «کارهای‌خانه برای فاطمه (س) و کارهای بیرون از خانه برای علی (ع)».

لیخند بر لبان فاطمه (س) نشست و گفت: «خدمای‌داند که من چقدر از این تقسیم خوشحالم»!

کارهای خانه کم نبود؛ اما زهرا (س) خوشحال و راضی بود. می‌گفت: «از سعادت‌مندی زن این است که بی‌دلیل در گذر نگاه نامحرم‌ها نباشد».

آرد کردن جو یا گندم تا پخت نان، طبخ و آماده کردن غذا همه بر دوش فاطمه (س) بود. کارهای مربوط به خانه یک طرف، رسیدگی به بچه‌ها و هم‌بازی شدن با آنها هم طرف دیگر.

علی (ع) آب خوردن تهیه می‌کرد، برای منزل هیزم می‌آورد، خریدهای خانه را انجام می‌داد و غیره؛ اما فقط کارهای بیرون منزل را انجام نمی‌داد. اگر فرصتی می‌یافت، خانه را جارو و به زهرایش در کارهای خانه کمک می‌کرد.

آن روز پیامبر (ص) مهمان خانه‌شان بود. مهمان زودتر از موعد رسید. دختر و دامادش را دید که با هم نشسته‌اند به پاک کردن عدس. با لیخندی که نشان از خرسندی و خوشحالی بود، گفت: «خدا به مردی که در کار خانه به همسرش کمک می‌کند به اندازه موهای بدنش ثواب عبادت می‌دهد».

('جامعه‌شناسی، 'کلیات اسلام'

کتاب «ادیان زنده جهان» اثر دین‌شناس فقید انگلیسی، رابرت. ا. هیوم یکی از منابع درسی و آکادمیک شناخته شده رشته‌های ادیان و عرفان و الهیات است. زمانی که این اثر برای اولین بار به بازار آمد، کتاب روزآمد دیگری در این زمینه در دسترس نبود؛ به همین دلیل این اثر خیلی زود جای خود را در میان استادان و دانشجویان رشته ادیان باز کرد.

هدف این اثر، بررسی دقیق منشاء متون مقدس، تحولات تاریخی و ارزش‌های اساسی ادیان جهان است. یازده دین در جهان وجود دارند که تنها دوتای آنها از مسیحیت جدیدترند. البته مذاهب جدید و نوظهوری هم در جهان فعال هستند که اکنون به صورت بین المللی درآمده‌اند و ادیانی هم از قدیم بوده‌اند که هیچ تأثیری در تمدن جهان نداشته‌اند. در این کتاب، تنها به یازده دینی پرداخته می‌شود که در طول تاریخ بشری به وجود آمده و پیوسته در طول اعصار، به حیات خود در سازمان مذهبی اجتماعی، هنر و ادبیات، و در قالب آداب دینی ادامه داده‌اند.

کتاب سیزده فصل دارد. دین هندو، دین جاپینی، بودایی، سیک، کنفوسیوسی، تائویی، شینتو، یهودی، زرتشتی، اسلام و مسیحیت به ترتیب عناوین فصول دوم تا دوازدهم کتاب‌اند. از جمله مطالب مهمی که درباره این ادیان در کتاب مطرح شده است. جایگاه‌شان در میان ادیان دیگر، سرگذشت بنیان‌گذاران، کتب مقدس و اصول ادیان یاد شده و نقاط ضعف و قوت آنها است.

فصل اول کتاب به مباحث نظری درباره دین می‌پردازد و فصل انتهایی نیز مقایسه‌ای اجمالی میان ادیان زنده جهان دارد و نکاتی کلی را درباره شباهت‌ها و تضادهای ادیان با یکدیگر، با تمرکز بر مسیحیت بیان می‌کند.

('جامعه‌شناسی، 'مدیریت و کسب و کار'

«معمای هابرماس» نوشته لاسه توماسن، استاد علوم سیاسی و روابط بی‌الملل دانشگاه اسکس است. این کتاب به بررسی آرا و اندیشه‌های یورگن هابرماس (۱۹۲۹-)، فیلسوف آلمانی در حوزه اقتصاد و سیاست می‌پردازد.

این کتاب درباره قانون و دموکراسی است و علی‌الخصوص به این موضوع می‌پردازد که چه عاملی به قانون مشروعیت می‌بخشد. به عبارت دیگر، دغدغه کتاب این سوال است که پیروی از قانون برای خود قانون است یا به سبب ترس از تلافی ناشی از قانون‌شکنی؟ این کتاب و این نقل قول، به مباحث طولانی‌مدت برابری و آزادی در

فلسفه و نظریه قانونی و سیاسی اشاره می‌کند: چه زمانی می‌توان گفت که شهروندان آزادند و با برابری از این آزادی لذت می‌برند؟ آزاد بودن به چه معناست؟ ... این کتاب مهم‌ترین موضوعات مطرح‌شده در آثار هابرماس و نیز مهم‌ترین جنبه‌های آن موضوعات را در بر می‌گیرد. در فصل اول، دیدگاه هابرماس در برابر سایر نظریه‌پردازان انتقادی بررسی می‌شود. در فصل دوم، به معرفی نوشته‌های هابرماس در باب حوزه عمومی، علی‌الخصوص اثر تاثیرگذارش، «دگرگونی ساختاری حوزه عمومی» می‌پردازیم.

در فصل سوم بر نظریه کنش ارتباطی تمرکز شده است. در این فصل، ایده‌های اصلی این نظریه را توضیح می‌دهم و شرح می‌دهم که چگونه هابرماس به این اندیشه‌ها رسید.

هابرماس بر پایه نظریه عقلانیت ارتباطی، اخلاق گفتمان (یا تا حدودی نظریه گفتمانی اعتبار) را مطرح می‌کند. این موضوع را در فصل چهارم بررسی می‌کنم. در این فصل نظریه اعتبار هابرماس، از جمله اعتبار هنجارهای اجتماعی تبیین می‌شود.

فصل پنج بر این موضوع متمرکز است که قانون هنگامی مشروعیت می‌یابد که مخاطبان آن همان مولفانش باشند.

در فصل ششم به بررسی سه موضوعی می‌پردازیم که از اواسط دهه ۱۹۹۰ در مرکز توجه هابرماس بوده است. نخستین مسئله به دولت ملت مربوط می‌شود. موضوع دوم، نقش مذهب در جوامع معاصر است. نهایتاً موضوع سوم به چالش‌های ناشی از تکنولوژی‌های جدید ژنتیکی برای فهم ما از انسان خودمختار می‌پردازد.