

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



دانشکده مهندسی برق و کامپیوتر

بازیابی هوشمند اطلاعات - تمرین پنجم

سید مهدی رضوی

استاد : خانم دکتر شاکری

دی ماه ۱۴۰۲



فهرست مطالب

۳	۱ تمرین اول
۶	۲ تمرین دوم
۹	۳ تمرین سوم

فهرست تصاویر

۳	Logistic Regression Results	۱
۴	Naive Bayes Results	۲
۵	SVM Results	۳
۶	Coherence Scores for topic modeling	۴
۷	Topics for 6 topics (1)	۵
۸	Topics for 6 topics (2)	۶

۱ تمرین اول

نتایج حاصل از اجرای سه مدل یادگیری ماشین ذکرشده در تمرین به شرح زیر خواهد بود. همانطور که از روی متن بالای هر نتیجه مشخص است، ما به طور کلی ۳ ویژگی بسیار معروف را برای این ۳ مدل یادگیری ماشین در نظر گرفته ایم. تعداد تکرار هر ترم، سپس TF-IDF و در نهایت هم از ویژگی همسایگی بین کلمات یا همان Ngram استفاده کرده ایم.

```
n_iter_i = _check_optimize_result(
----- RESULTS FOR NGRAM FEATURE-----
Validation Accuracy (Logistic Regression)(NGRAM) : 0.7062825130052021
Validation Precision (Logistic Regression)(NGRAM) : 0.6758266671056454
Validation Recall (Logistic Regression)(NGRAM) : 0.7062825130052021
Validation F1-Score (Logistic Regression)(NGRAM) : 0.6878520149969166
-----

----- RESULTS FOR TFIDF FEATURE-----

Validation Accuracy (Logistic Regression)(TFIDF) : 0.7062825130052021
Validation Precision (Logistic Regression)(TFIDF) : 0.6758266671056454
Validation Recall (Logistic Regression)(TFIDF) : 0.7062825130052021
Validation F1-Score (Logistic Regression)(TFIDF) : 0.6878520149969166
-----

Validation Accuracy (Logistic Regression): 0.7078831532613046
Test Accuracy (Logistic Regression): 0.7130852340936374
----- RESULTS FOR COMBINED FEATURES -----
Validation Accuracy (Logistic Regression) : 0.7078831532613046
Validation Precision (Logistic Regression) : 0.6791136388719411
Validation Recall (Logistic Regression): 0.7078831532613046
Validation F1-Score (Logistic Regression) : 0.6903080976263811
```

شکل ۱: Logistic Regression Results

به نظر می رسد که تفاوت قابل توجهی در عملکرد بین استفاده از ویژگی n-gram و ویژگی TF-IDF به صورت جداگانه وجود ندارد. با این حال، هنگام ترکیب هر دو ویژگی، بهبود جزئی در برخی از معیارها وجود دارد. در اینجا چند مشاهدات بر اساس نتایج آورده شده است:

دقت: دقت برای هر دو ویژگی n-gram و ویژگی TF-IDF به صورت جداگانه در (0.7063) یکسان است.

در صورت ترکیب، دقت کمی بهبود می یابد و به (0.7079) می رسد. این نشان می دهد که ترکیب ویژگی ها ممکن است تأثیر مثبت کوچکی بر دقت کلی داشته باشد.

دقت، یادآوری و امتیاز: F1 مقادیر دقت، فراخوانی و امتیاز F1 برای هر دو ویژگی فردی یکسان است.

با این حال، هنگام ترکیب ویژگی ها، بهبود جزئی در دقت (0.6791) و امتیاز F1 (0.6903) مشاهده می شود. این نشان می دهد که ترکیب ویژگی ها ممکن است به دقت بهتر و عملکرد کلی مدل کمک کند.

بر اساس این نتایج، به نظر می رسد که ترکیبی از هر دو ویژگی n-gram و TF-IDF در مقایسه با استفاده از هر ویژگی به صورت جداگانه، بهبودی حاشیه ای در عملکرد مدل ارائه می دهد. در حالی که این بهبود قابل توجه نیست، نتیجه گیری کلی ما از این آزمایش این خواهد بود که: اگر به دنبال پیشرفت های افزایشی در دقت و دقت کلی مدل هستید، ارزش آن را دارد که رویکرد ویژگی ترکیبی را در نظر بگیرید.

```
-----
Validation Accuracy (Naive Bayes)(TF-IDF) : 0.720688275310124
Test Accuracy (Naive Bayes)(TF-IDF) : 0.7078831532613046
-----
----- RESULTING FOR TF-IDF FEATURE -----
Validation Accuracy (Naive Bayes)(TF-IDF) : 0.720688275310124
Validation Precision (Naive Bayes)(TF-IDF) : 0.6113871059869905
Validation Recall (Naive Bayes)(TF-IDF) : 0.720688275310124
Validation F1-Score (Naive Bayes)(TF-IDF) : 0.661237640343044
-----
Validation Accuracy (Naive Bayes)(NGRAM) : 0.7242897158863545
Test Accuracy (Naive Bayes)(NGRAM) : 0.7158863545418167
-----
----- RESULTING FOR NGRAM FEATURE -----
Validation Accuracy (Naive Bayes)(NGRAM) : 0.7242897158863545
Validation Precision (Naive Bayes)(NGRAM) : 0.6380881769695498
Validation Recall (Naive Bayes)(NGRAM) : 0.7242897158863545
Validation F1-Score (Naive Bayes)(NGRAM) : 0.6672986783818469
-----
Validation Accuracy (Naive Bayes)(Counting) : 0.7130852340936374
Test Accuracy (Naive Bayes)(Counting) : 0.7090836334533813
-----
----- RESULTS FOR COUNTING FEATURE -----
Validation Accuracy (Naive Bayes)(Counting) : 0.7130852340936374
Validation Precision (Naive Bayes)(Counting) : 0.6643255132793805
Validation Recall (Naive Bayes)(Counting) : 0.7130852340936374
Validation F1-Score (Naive Bayes)(Counting) : 0.6742988238734677
```

شکل ۲: Naive Bayes Results

از این نتایج می توان موارد زیر را مشاهده کرد: ویژگی NGRAM به بالاترین دقت اعتبار سنجی دست می یابد و به دنبال آن ویژگی TF-IDF قرار دارد. با این حال، هنگام در نظر گرفتن دقت تست، ویژگی TF-IDF کمی بهتر عمل می کند. دقت، فراخوانی و امتیاز F1 الگوهای مشابهی را در بین ویژگی ها نشان می دهد، با ویژگی NGRAM به طور کلی بهتر از ویژگی های TF-IDF و شمارش. شایان ذکر است که معیارهای عملکرد ممکن است بسته به مجموعه داده خاص و ماهیت مشکل طبقه بندی متفاوت باشد. بر اساس این نتایج، به نظر می رسد که ویژگی NGRAM عملکرد کمی بهتر از خود نشان می دهد و در مقایسه با ویژگی های TF-IDF و شمارش، به دقت و امتیاز F1 بالاتری دست می یابد.

```
Validation Accuracy (SVM) : 0.7222889155662265
Validation Precision (SVM) : 0.6527118658402672
Validation Recall (SVM) : 0.7222889155662265
Validation F1-Score (SVM) : 0.672308578987064
*****
----- RESULTS FOR COUNTING FEATURE -----
Validation Accuracy (SVM) : 0.7014805922368947
Validation Precision (SVM) : 0.6423272638890312
Validation Recall (SVM) : 0.7014805922368947
Validation F1-Score (SVM) : 0.6489195413892614
*****
----- RESULTS FOR SVM NGRAM FEATURE -----
Validation Accuracy (SVM)(NGRAM) : 0.7078831532613046
Validation Precision (SVM)(NGRAM) : 0.6577867787101905
Validation Recall (SVM)(NGRAM) : 0.7078831532613046
Validation F1-Score (SVM)(NGRAM) : 0.6539472765302519
*****
----- RESULTS FOR TF-IDF FEATURE -----
Validation Accuracy (SVM) : 0.7222889155662265
Validation Precision (SVM) : 0.6527118658402672
Validation Recall (SVM) : 0.7222889155662265
Validation F1-Score (SVM) : 0.672308578987064
```

شکل ۳: SVM Results

از این نتایج می توان موارد زیر را مشاهده کرد:

مدل SVM با ویژگی شمارش کمترین عملکرد را در تمام معیارها، با دقت، دقت، فراخوانی و امتیاز F1 کمتر در مقایسه با سایر ویژگی ها به دست می آورد. هر دو ویژگی NGRAM و TF-IDF عملکرد مشابهی را نشان می دهند و در مقایسه با ویژگی شمارش، به دقت، دقت، فراخوانی و امتیاز F1 بالاتری دست می یابند. مدل SVM با ویژگی TF-IDF به بالاترین دقت اعتبارسنجی و امتیاز F1 دست می یابد، در حالی که ویژگی NGRAM دقت و یادآوری کمی بالاتری را نشان می دهد. بر اساس این نتایج، به نظر می رسد که مدل SVM با ویژگی TF-IDF به طور کلی بهترین عملکرد را دارد و دقت و امتیاز F1 بالاتری را در مقایسه با سایر ویژگی ها نشان می دهد.

در نهایت می توان به این تبصره کلی اشاره کرد که این نتایج به دست آمده مربوط به داده های ما می باشد و ممکن است با توجه به تغییر بعضی هایپرپارامترها بتوان تغییرات گسترده تری در آن ها به وجود آورد.

۲ تمرین دوم

نتایج مربوط به آزمایش ما برای مدل کردن Topic ها به شرح زیر است. به ازای هر کدام از تعداد Topic بین ۳ تا ۱۵ امتیاز Coherent Score به شرح زیر است :

Coherence Score for the 3 topics : 0.6767713546654176
Coherence Score for the 4 topics : 0.6849739873787832
Coherence Score for the 5 topics : 0.6953981779587295
Coherence Score for the 6 topics : 0.6983948244835263
Coherence Score for the 7 topics : 0.6298256798505133
Coherence Score for the 8 topics : 0.5388533596855748
Coherence Score for the 9 topics : 0.6358683188913125
Coherence Score for the 10 topics : 0.4694451804427812
Coherence Score for the 11 topics : 0.4911213140781505
Coherence Score for the 12 topics : 0.5546870641075222
Coherence Score for the 13 topics : 0.45232996654013785
Coherence Score for the 14 topics : 0.521629690579584
Coherence Score for the 15 topics : 0.5369786191412719

شکل ۴: Coherence Scores for topic modeling

(0, 'علم و The Verge: علم و تکنولوژی ۳" + "0.004* منبع: GSM Arena منبع: "0.005*
'علم و تکنولوژی ۳" + "0.001* برای آگاهی: gsmarena: تکنولوژی ۳" + "0.002* منبع:
'از آخرین اخبار و اطلاعات جشنواره فیلم فجر به صفحه ویژه جشنواره فیلم فجر ۹۷
' + "علم و تکنولوژی ۳ Sam Mobile: کالا مگ بروید." + "0.001* منبع: u200cدر دیجی
'خرید کتاب از ۳۰ off: سلامت و زیبایی ۲" + "0.001* کد تخفیف: BBC منبع: "0.001*
'بازی ویدیویی ۰." + "0.000* این VG ۲۴۷: فیدیبو کتاب و ادبیات ۶" + "0.001* منبع:
'توانید آن را با ۳۰ درصد تخفیف از u200c\کتاب در کمپ-ن-جی بخونم قرار دارد و می
'علم و تکنولوژی ۱ MOTOR: سایت فیدیبو دانلود و مطالعه کنید." + "0.000* منبع:
'"), (۳)
(1, 'بازی Polygon: علم و تکنولوژی ۳" + "0.003* منبع: Phone Arena منبع: "0.004*
'VG ۲۴۷: بازی ویدیویی ۰." + "0.003* منبع: Gematsu: ویدیویی ۰." + "0.003* منبع:
'توانید با ۳۰ درصد تخفیف از u200c\بازی ویدیویی ۰." + "0.002* این کتاب را می
'سایت فیدیبو دانلود و مطالعه کنید." + "0.002* منبع: مارک براون - یوتیوب بازی
' + "ویدیویی ۰." + "0.001* بی." + "0.001* دانلود کتاب از فیدیبو کتاب و ادبیات ۶
'"), (0.001* منبع: To a Mac جی: "۹" + "0.001* منبع: "۹")
(2, 'علم و تکنولوژی ۳ Android Authority: سلامت و زیبایی ۲" + "0.004* منبع: "0.025*
'The Verge: علم و تکنولوژی ۳" + "0.001* منبع: GSMArena منبع: "0.004* +
'ها با ما و دیگران در میان u200c\سینما ۵" + "0.001* نظرات خود را در بخش کامنت
'سلامت و زیبایی ۲" + "0.001* شما در این Engadget: بگذارید." + "0.001* منبع:
'کنید؟" + "0.001* خیلی اوقات از دور تماشا کردن و عکس u200c\زمینه چه فکر می
'لطفاً به خصوص وقتی جاده شلوغ است، درست رانندگی u200c\گرفتن لذت بیشتری دارد
'کنید لازم نیست همیشه تخت گاز برانیم، لازم نیست دائماً در تلاش برای سبقت
'گرفتن از ماشین جلویی باشیم و لازم نیست هرسال اثبات کنیم که بعد از رانندگان
'تانزانیا، لیبی، مالای و جمهوری دموکراتیک کنگو، بدترین رانندگان جهان
' "هستیم!" + "0.001* (اصلاح این یکی که دیگر دست خودمان است.) سلامت و زیبایی ۲
'دستور پخت و عکس: مریم ابراهیمی بیشتر بخوانید: چطور کیک سب و دارچین" + "0.001* +
'"), (درست کنیم؟)
(3, 'The Verge: هنر و سینما ۵" + "0.011* علم و تکنولوژی ۳" + "0.003* منبع: "0.023*
'علم و تکنولوژی ۳" + "0.003* راهنمای خرید ۱" + "0.003* برای مطالعه روی عکس زیر
'بازی ویدیویی ۰." + "0.002* اگر به این The Verge: کلیک کنید." + "0.002* منبع:
'مطلب علاقه داشتید، سایر مطالب مجله را هم مطالعه کنید." + "0.001* منتقدان چه
'علم و تکنولوژی ۳" + "0.001* در این motor: گویند؟" + "0.001* منبع: u200c\می
'ها و شعرهای به u200c\های خوب، قصه u200c\برنامه سعی کردیم موسیقی u200c\ویژه
'یادماندنی را برایتان انتخاب کنیم و همچنین هر روز یکی از بهترین
'، های سینمایی ایران و جهان را به انتخاب مازیار وکیلی منتقد سینما u200c\فیلم

Topics for 6 topics (1) : ۵ شكا

'گرفتن از ماشین جلویی باشیم و لازم نیست هرسال اثبات کنیم که بعد از رانندگان
'تانزانیا، لیبی، مالاوی و جمهوری دموکراتیک کنگو، بدترین رانندگان جهان
'هستیم!" + 0.001*0.001)"(اصلاح این یکی که دیگر دست خودمان است...) سلامت و زیبایی ۲
'دستور پخت و عکس: مریم ابراهیمی بیشتر بخوانید: چطور کیک سیب و دارچین" + 0.001*0.001
'("درست کنیم؟")

(3,
'TheVerge :هنر و سینما ۵" + 0.011*0.003 + 0.003*0.003 "منبع" *0.023
'علم و تکنولوژی ۳" + 0.003*0.003 "راهنمای خرید ۱" + 0.003*0.003 "برای مطالعه روی عکس زیر
'بازی ویدیویی ۰" + 0.002*0.002 "اگر به این The Verge :کلیک کنید." + 0.002*0.002 "منبع"
'مطلب علاقه داشتید، سایر مطالب مجله را هم مطالعه کنید." + 0.001*0.001 "منتقدان چه
'علم و تکنولوژی ۳" + 0.001*0.001 "در این ۱ motor :گویند؟" + 0.001*0.001 "منبع" *0.001
'ها و شعرهای به \u200c\u200cهای خوب، قصه \u200c\u200cبرنامه سعی کردیم موسیقی \u200c\u200cویژه
'یادماندنیی را برایتان انتخاب کنیم و همچنین هر روز یکی از بهترین
'،های سینمایی ایران و جهان را به انتخاب مازیار وکیلی منتقد سینما \u200c\u200cفیلم
'("کنیم \u200c\u200cپیشنهاد می
(4,
' + "بازی ویدیویی . GameSpot :بازی ویدیویی ۰" + 0.002*0.002 "منبع" *0.009
'سلامت و زیبایی Inc: علم و تکنولوژی ۳" + 0.001*0.001 "منبع" phonearena "منبع" *0.001
' + "ارزش ۳۰۰ دلار را دارد؟ Nintendo Switch قسمت شانزدهم: آیا" + 0.001*0.001 + ۲
'های پاییزی را از دست ندهید (قسمت اول) قسمت \u200c\u200cقسمت چهارم: این بازی" *0.001
'های پاییزی را از دست ندهید (قسمت دوم) قسمت ششم: جادوی \u200c\u200cپنجم: این بازی
'بخیریم یا نه؟" + 0.001*0.001 "قسمت نوزدهم: آخرین Nintendo Switch :فوتبال قسمت هفتم
'مخاف سونی تو زرد از آب درآمد!" + 0.001*0.001 "قسمت دوازدهم: پایان ده سال انتظار:
'چرا فاینال فانتزی ۱۵ را دوست داریم؟" + 0.001*0.001 "قسمت سیزدهم: رزیدنت اوایل ۷
'از گور برآمد هفتم قسمت چهاردهم: زمستان امسال چه بازی کنیم؟" + 0.001*0.001 "قسمت هفدهم
'های جدید مارول؛ مروری بر داستان کونا قسمت \u200c\u200cاز رزیدنت اوایل تا بازی
'("هجدهم: آیا رزیدنت اوایل ۷ طرفداران را راضی کرده است؟")

(5,
'سلامت و زیبایی ۲" + 0.001*0.001 "آخرین اخبار و اطلاعات Healthline :منبع" *0.001
'ی \u200c\u200cویژه \u200c\u200cهای مختلف جشنواره را در صفحه \u200c\u200cفیلم \u200c\u200cدرباره
'جشنواره فیلم فجر ۹۸ بخوانید." + 0.001*0.001 "بیشتر بخوانید: در روزهای
'"ی دیجی کالا مگ بشوید هنر و سینما \u200c\u200cفیلم فجر ۹۸ نویسنده \u200c\u200cجشنواره
'ها عرضه \u200c\u200cاستیشن ۵ به این زودی \u200c\u200cقسمت پنجاه و چهارم: آیا پلی" *0.001 +
' :شود؟" + 0.001*0.001 "منبع: یوروگیمز بازی ویدیویی ۰" + 0.001*0.001 "منبع" *0.001
' + "سلامت و زیبایی ۲ The Verge :علم و تکنولوژی ۳" + 0.001*0.001 "منبع" theverge
' + "هنر و سینما ۵" + 0.001*0.001 "کتاب و ادبیات ۶ Polygon :منبع" *0.001
'بیشتر بخوانید: نقد روزی روزگاری در هالیوود؛ وقتی مخدر سینما آرام زیر" *0.000
'("خزد هنر و سینما \u200c\u200cپوست می

شکل ۶: (2) Topics for 6 topics

۳ تمرین سوم

$$P(\text{Minus}|\text{Doc1}) < P(\text{Plus}|\text{Doc1})$$

$$P(\text{Minus}|\text{Doc2}) > P(\text{Plus}|\text{Doc2})$$

$$P(\text{Minus}|\text{Doc3}) < P(\text{Plus}|\text{Doc3})$$

$$P(\text{Minus}|\text{Doc4}) > P(\text{Plus}|\text{Doc4})$$

$$P(\text{Minus}|\text{Doc5}) < P(\text{Plus}|\text{Doc5})$$

با توجه به نتایج بالا ، هر یک از اسناد برجسب متناظر خود را دریافت خواهند کرد.
بر اساس منطق این رده‌بند ، هر کدام که از این اسناد که احتمال شرطی بالاتری برای هر کدام از کلاس‌ها داشته‌باشد ، به آن کلاس مربوط می‌شود.
باید به مقایسه دو عبارت زیر به ازای سند ۶ و ۷ بپردازیم و سپس برجسب متناظر با آن داده را بر روی سند مدنظر بزنیم.

$$P(\text{plus}) * P(W_1|\text{plus}) * P(W_2|\text{plus}) * P(W_3|\text{plus}) * \dots * P(W_n|\text{plus})$$

$$P(\text{minus}) * P(W_1|\text{minus}) * P(W_2|\text{minus}) * P(W_3|\text{minus}) * \dots * P(W_n|\text{minus})$$

با توجه به معادلات جایگاه نامعادلات بالا ، به داده‌های زیر خواهیم رسید :

$$P(\text{plus}) = \frac{3}{5}$$

$$P(\text{minus}) = \frac{2}{5}$$

$$P(\textit{love}|\textit{plus}) = \frac{1+1}{11+6} = \frac{2}{17}$$

$$P(\textit{movie}|\textit{plus}) = \frac{4+1}{11+6} = \frac{5}{17}$$

$$P(\textit{great}|\textit{plus}) = \frac{2+1}{11+6} = \frac{3}{17}$$

$$P(\textit{good}|\textit{plus}) = \frac{2+1}{11+6} = \frac{3}{17}$$

$$P(\textit{acting}|\textit{plus}) = \frac{1+1}{11+6} = \frac{2}{17}$$

$$P(I|\textit{plus}) = \frac{1+1}{11+6} = \frac{2}{17}$$

$$P(\textit{hated}|\textit{minus}) = \frac{1+1}{5+5} = \frac{2}{10}$$

$$P(I|\textit{minus}) = \frac{1+1}{5+5} = \frac{2}{10}$$

$$P(\textit{movie}|\textit{minus}) = \frac{1+1}{5+5} = \frac{2}{10}$$

$$P(\textit{poor}|\textit{minus}) = \frac{1+1}{5+5} = \frac{2}{10}$$

$$P(\textit{acting}|\textit{minus}) = \frac{1+1}{5+5} = \frac{2}{10}$$

Document6 : I loved the poor play.

$$P(plus|Document6) = P(plus) * \frac{P(love|plus) * P(I|plus) * P(poor|plus) * P(play|plus)}{P(Document6)}$$

$$P(minus|Document6) = P(minus) * \frac{P(love|minus) * P(I|minus) * P(poor|minus) * P(play|minus)}{P(Document6)}$$

با توجه به این که مقدار احتمال سند ششم برای هر دو ترم با هم برابر است ، به صورت زیر تخمین خواهیم زد :

$$P(plus|Document6) = \frac{3}{5} * \frac{2}{17} * \frac{2}{17} * \frac{1}{17} * \frac{1}{17} = 0.000028735$$

$$P(minus|Document6) = \frac{2}{5} * \frac{1}{10} * \frac{2}{10} * \frac{2}{10} * \frac{1}{10} = 0.00016$$

برچسب سند ششم منفی خواهد بود.

Document7 : I hated the play movie.

همانند جایگاه معادلات بالا برای سند ۶ ، مقادیر احتمالاتی را برای کلاس‌های ممکن برای سند ۷ را محاسبه خواهیم کرد :

$$P(plus|Document7) = \frac{3}{5} * \frac{2}{17} * \frac{1}{17} * \frac{1}{17} * \frac{5}{17} = 0.00007183822$$

$$P(minus|Document7) = \frac{2}{5} * \frac{2}{10} * \frac{2}{10} * \frac{2}{10} * \frac{1}{10} = 0.00032$$

برچسب سند هفتم منفی خواهد بود.

با توجه به نتایج بالا و وزن کلمه مثبتی همچون love و وزن منفی کلمه‌ای چون poor می‌توان پیش‌بینی نمود که وزن بیشتر به سمت کلاس منفی باشد ، چرا که احتمال منفی به شرط کلمه poor بسیار بیشتر از احتمال مثبت به شرط کلمه love است.

$$P(minus|poor) > P(plus|love)$$

و همچنین وزن کلمه منفی مانند hate همچنین نتایجی قابل پیش‌بینی بود.

نتایج به دست‌آمده تا حد خوبی منطقی هستند.

این نتایج به نظر بیشتر وابسته به داده‌های پس‌زمینه Background ما هستند. وابستگی میزان احتمال هر کلاس به شرط هر کلمه که از داده‌های پس‌زمینه ما به دست آمدند این قضیه را بیشتر اثبات می‌کند.