



بازیابی هوشمند اطلاعات

تمرین چهارم

نام و نام خانوادگی : حسین سیفی

شماره دانشجویی : ۸۱۰۱۰۰۳۸۶

دی ۱۴۰۰

## فهرست

بخش ۱ – Word Association ..... ۳

سوال ۱ ..... ۳

سوال ۲ ..... ۳

سوال ۳ ..... ۴

سوال ۴ ..... ۴

پیاده سازی ..... ۴

بخش ۲ – Clustering ..... ۵

سوال ۱ ..... ۵

سوال ۲ ..... ۶

سوال ۳ ..... ۶

پیاده سازی ..... ۶

## بخش ۱ – Word Association

در این تمرین هدف استخراج روابط syntagmatic و paradigmatic بین کلمات با استفاده از Mutual Information است. روش کار به این صورت است که مقدار MI را برای هر زوج کلمه محاسبه می‌کنیم. مقادیر MI برای یک کلمه در مواجهه با کلمات دیگر نشان دهنده روابط Paradigmatic است و با ضرب همین بردار به عنوان بردار Context در بردار Context کلمات دیگر، می‌توان روابط Syntagmatic را به دست آورد.

به عبارتی دیگر با به دست آوردن بردار MI هر کلمه در صورت مواجهه با کلمات دیگر نشان می‌دهیم که کدام کلمات بیشتر با همنشینی یکدیگر ظاهر شده‌اند تا به صورت تکی و با محاسبه شباهت کسینوسی همین مقادیر به عنوان بردار Context در بردار Context کلمات دیگر، بررسی می‌کنیم که کدام کلمات با کلمات مشابهی همنشین هستند و در نتیجه قابل جانشینی با یکدیگر هستند.

### سوال ۱

برای به دست آوردن بردار Context هر کدام از کلمات از فرمول زیر استفاده می‌کنیم.

$$I(X_{w1}; X_{w2}) = \sum_{u \in \{0, 1\}} \sum_{v \in \{0, 1\}} p(X_{w1} = u, X_{w2} = v) \log_2 \frac{p(X_{w1} = u, X_{w2} = v)}{p(X_{w1} = u)p(X_{w2} = v)}$$

با استفاده از فرمول فوق و کمک گرفتن از ماتریس رویداد همزمان کلمات، می‌توان Mutual Information را برای هر زوج کلمه در لغت‌نامه (Vocabulary) اسناد محاسبه کرد و در ادامه سوالات بخش ۱ مقدار MI یک کلمه به صورت زوج با سایر کلمات به عنوان بردار Context و برای یافتن روابط همنشینی و جانشینی کاربرد دارد.

### سوال ۲

ده کلمه‌ای که بیشترین روابط Syntagmatic را با کلمات Teacher و Iran دارند، به شرح زیر می‌باشند:

<i>Iran</i>	<i>Teacher</i>	رتبه
algier	fleischer	۱
battlefield	hardli	۲
battleship	strap	۳
civilian	brasilia	۴
colder	immun	۵
dissoci	laden	۶
efficaci	accord	۷
enemi	adjust	۸
es	associ	۹
homeland	began	۱۰

### سوال ۳

ده کلمه‌ای که بیشترین روابط Paradigmatic را با کلمات Teacher و Iran دارند ، به شرح زیر می‌باشند:

<i>Iran</i>	<i>Teacher</i>	رتبه
Iraq	Laden	۱
Iranian	Brasilia	۲
Missil	Grind	۳
Tehran	Strap	۴
Hormuz	Hardli	۵
Gulf	Fleischer	۶
Attack	Coalit	۷
Naval	Unrest	۸
Warship	Junior	۹
Chines	Bleak	۱۰

### سوال ۴

با بررسی جداول بالا متوجه می‌شویم که هر کدام از کلمات همنشین با کلمه Iran به نوعی با این کلمه مرتبط هستند و در اخبار متفاوت می‌توانند همزمان با کلمه Iran در متون متفاوت ظاهر شوند. همچنین با بررسی بردار Context کلمات ذکر شده که دارای مقدار بالایی برای رابطه Syntagmatic با کلمه Iran هستند مشخص می‌شود که این کلمات بردار Context تقریباً مشابه دارند و در نتیجه با کلمات یکسانی رابطه Paradigmatic دارند پس این کلمات می‌توانند جایگزین کلمه Iran باشند.

از طرفی دیگر به نظر می‌رسد کلمات دارای روابط Syntagmatic با کلمه Teacher ارتباط معنایی چندانی با آن ندارند در نتیجه با توجه به اینکه همین روابط Syntagmatic به عنوان بردار Context این کلمه استفاده می‌شوند ، روابط Paradigmatic نیز اعتبار پایینی به دست می‌آورند. با بررسی کلمات دارای ارتباط paradigmatic با کلمه Teacher متوجه می‌شویم که حدس درستی داشته‌ایم و به نظر نمی‌رسد که بتوانیم کلمه Teacher را با این کلمات جایگزین کنیم. دلیل این رویداد این است که کلمه Teacher احتمال بسیار پایینی برای ظهور در اخبار (برابر ۰.۰۰۰۲۷) دارد و تنها در یک سند ظاهر شده است.

### پیاده سازی

برای پیاده سازی این تمرین از زبان پایتون و کتابخانه‌های متعددی استفاده شده است که به صورت خلاصه توضیح داده می‌شود. در این تمرین از کتابخانه nltk برای تبدیل اخبار به بردارهایی از کلمات و تکرار هر کدام استفاده شده است و از کتابخانه‌های دیگری مانند pandas ، numpy ، scikit-learn ، scipy ، tqdm و enchant کمک گرفتیم. در ابتدای کار فایل csv اخبار به کمک کتابخانه pandas داخل برنامه بارگزاری می‌شود و پیش‌پردازشی شامل حذف علائم نگارشی ، حذف کلمات بی‌معنی ، حذف اعداد و ریشه‌یابی روی اخبار به کمک کتابخانه‌های مرتبط انجام می‌گیرد. در ادامه تعداد تکرار کلمات در سندهای مختلف محاسبه می‌شود ، لغت‌نامه‌ای شامل تمام کلمات اخبار ایجاد می‌شود ، تعداد سندهای شامل هر کلمه مشخص می‌شود ، کلمات پرتکرار و عمومی که تکراری بیشتر از ۴۰۰ بار دارند حذف می‌شوند و ماتریسی شامل تعداد ظهور همزمان هر جفت کلمه ساخته می‌شود. سپس مقدار

mutual Information برای هر جفت کلمه با استفاده از فرمول ذکر شده در سوال ۱ محاسبه می‌شود. این مقادیر نشان دهنده ارتباط همنشینی می‌باشند و با محاسبه فاصله کسینوسی بین بردار کلمات همنشین (Context) کلمات مختلف می‌توان روابط جانشینی را برای هر کلمه به دست آورد.

## بخش ۲ – Clustering

در بخش ۲ این تمرین به Train کردن مدل‌های متفاوت خوشه بندی با استفاده از Dataset اخبار و ارزیابی بر اساس Label‌های موجود برای هر کدام از این اخبار می‌پردازیم.

### سوال ۱

مقدار چهار معیار ارزیابی خواسته شده برای روش kmeans به شرح زیر می‌باشد:

مقدار	معیار/ارزیابی
۰.۷۹۷	Purity
۰.۳۹۰	F1
۰.۴۸۴	NMI
۰.۶۹۷	RI

برای محاسبه معیارهای F1 ، NMI و RI توابع آماده کتابخانه scikit-learn استفاده شد و برای محاسبه Purity با اندکی جستجو و کمک گرفتن از ماتریس contingency نتیجه مطلوب به دست آمد. همچنین تعداد False Positive و False Negative به تفکیک کلاس در جدول زیر آورده شده است:

False Positive	False Negative	کلاس
۵۳۸	۹۷	Acq
۱۸	۲۷	Crude
۵۴	۴۱۳	Earn
۰	۴۱	Grain
۰	۱۹۱	Interest
۰	۲۲۲	Money-fx
۰	۱۰۸	Ship
۵۰۳	۱۴	Trade

مشاهده می‌شود که بیشترین تعداد False Negative مربوط به کلاس Earn و بیشترین تعداد False Positive مربوط به کلاس Acq می‌باشد. متأسفانه کلاس‌های کوچکتر به دلیل اینکه در هیچکدام از کلاس‌های مقصد بیشترین تعداد را نسبت به دیگر کلاس‌ها نداشتند ، هیچ کلاس اخصیص یافته‌ای ندارند و در نتیجه تمامی تعداد آنها به عنوان False Negative معرفی می‌شود.

## سوال ۲

مقدار معیارهای ارزیابی ذکر شده برای سه روش سلسله مراتبی Single link ، Average link و Complete link به شرح زیر می باشد:

معیار ارزیابی	Complete link	Average link	Single link
Purity	۰.۶۱۴	۰.۵۵۳	۰.۵۱۷
F1	۰.۱۶۱	۰.۱۵۰	۰.۰۸۶
NMI	۰.۱۴۷	۰.۳۱۷	۰.۰۰۲
RI	۰.۶۲۴	۰.۶۳۳	۰.۳۵۸

به نوعی مهمترین معیار ارزیابی را می توان Normalized Mutual Information معرفی کرد و با توجه به جدول فوق بهترین مقدار برای این معیار در روش Average link و سپس با اختلاف کمی در روش Complete link به دست می آید اما معیاری مانند Purity اطلاعات کاملی به ما نمی دهد و فقط یکسان بودن اعضای خوشه های پیشبینی شده را به ما نشان می دهد.

## سوال ۳

در مقایسه ی این چهار روش می توان بگوییم که بین این روش ها بهترین مقدار Rand Index متعلق به روش kmeans می باشد و این معیار نشان می دهد که Kmeans بیشترین مقدار داده هایی را دارد که در دسته ی درست قرار گرفته اند. سایر معیارهای ارزیابی Purity ، NMI و F1 نیز بیشترین مقدار خود را در روش Kmeans دارند. این نتایج به این معنی است که روش kmeans بیشترین مقدار یکسان بودن خوشه ها و بیشترین مقدار تعادل بین Precision و Recall را دارد. پس در مجموع بهترین روش بین روش های بررسی شده Kmeans می باشد و Average link با اختلاف بسیار زیادی پشت آن ایستاده است.

## پیاده سازی

در پیاده سازی این بخش نیز همانند بخش ۱ ابتدا به پیش پردازش روی مجموعه داده می پردازیم و کلمات و کاراکترهایی که کمترین کمک را به ما در خوشه بندی اخبار می کنند یا ما را گمراه می کنند ، حذف می کنیم. ادامه کار از جمله به دست آوردن احتمال هر کلمه ، ساختن مدل ها و انجام تنظیمات آنها ، محاسبه معیارهای ارزیابی و به دست آوردن برچسب های پیشبینی شده توسط توابع پیشفرض کتابخانه scikit-learn انجام می گیرد.