

گزارش تکلیف Cosine Similarity

درس داده کاوی

امیرحسین ابوالحسنی

۴۰۰۴۰۵۰۰۳

مقدمه

برای ارزیابی شباهت (Similarity) و یا تفاوت (Disimilarity) دو بردار از معیارهای متفاوتی می توان استفاده نمود. یکی از این معیارها تحت عنوان Cosine Similarity شناخته می شود که از این رابطه محاسبه می گردد:

$$\cos \theta = \frac{\vec{u} \cdot \vec{v}}{||\vec{u}|| \cdot ||\vec{v}||}$$

که $\cos \theta$ همان Cosine Similarity می باشد.

حال با توجه به این متریک برای بررسی شباهت دو بردار، سعی در شباهت سنجی دو متن خبری فارسی می گردد. برای این کار سه خبر از خبرگزاری تسنیم تهیه شده است که موضوعات آنها طبق این جدول می باشد.

شماره خبر	نام خبر	موضوع
خبر ۱	آزمون: می دانیم چطور مقابل ازبکستان بازی کنیم	ورزشی - فوتبال - ایران
خبر ۲	جمع آوری ۳ "پتابایت" داده های ماهواره ای توسط سایت ورامین	فضا و نجوم
خبر ۳	بگیرستاین در پایان فصل منچستر سیتی را ترک می کند	ورزشی - فوتبال - جهان

در این تکلیف، هر کدام از این اخبار شباهتشان با یکدیگر سنجیده می شود. انتظار می رود دو خبر ورزشی بیشترین مقدار شباهت را داشته باشند.

کتابخانه ها

در این تکلیف برای بخش Web Scraping از کتابخانه های:

• Beautiful Soup

• Requests

و برای پردازش متن از کتابخانه های:

• hazm: پردازش متن فارسی

• nltk: شناسایی ترکیب های مورد علاقه

• re

و برای محاسبه Cosine Similarity از کتابخانه Numpy استفاده شده است.

Web Scraping

ابتدا با استفاده از کتابخانه Requests صفحه HTML گرفته شده، سپس با استفاده از کتابخانه BeautifulSoup این HTML پارس شده و متن خبر استخراج می‌گردد.

پیش پردازش متن

۱.۰ Normalization

ابتدا متن خام، نرمالایز^۱ می‌شود تا حروف و کلمات اضافی و غیرقابل استفاده از آن حذف شود.

۲.۰ Tokenization

در این مرحله، ابتدا متن را جمله جمله کرده، و سپس کلمه کلمه می‌کنیم تا تمامی کلمات موجود در متن را به شکل لیست در اختیار داشته باشیم.

۳.۰ برچسب گذاری دستوری

به هر کلمه یک برچسب تخصیص داده می‌شود که نشان دهنده نقش دستوری آن کلمه در آن جمله می‌باشد.

برای مثال:

اسم - فعل - قید - ...

بازیابی کلمات کلیدی

گرامرهایی تعریف می‌کنیم که عباراتی با این گرامر برای ما اهمیت دارند. برای مثال گرامر

NP: <NOUN.*><ADJ.*>? # Noun(s) + Adjective(optional)

به معنی عباراتی هستند که از یک اسم و یک صفت تشکیل شده اند.

Vectorization

به استفاده از مدل Sec2Vec دنباله کلمات کلیدی را به فضای برداری می‌بریم.

محاسبه شباهت کسینوسی

با توجه به فرمول ارائه شده و کتابخانه Numpy این مقدار را برای دو بردار محاسبه می‌کنیم.

نتیجه گیری

شباهت کسینوسی هر دو خبر به صورت زیر است:

خبر ۱	خبر ۲	خبر ۳	
۱	۰.۳۴	۰.۳۶	خبر ۱
	۱	۰.۲۲	خبر ۲
		۱	خبر ۳

همانطور که انتظار می‌رفت، خبر ۱ و ۳ بیشترین شباهت را نسبت به هر دو خبر دیگری دارند.

^۱Normalize