



دسته‌بندی متون پروژه مبانی برنامه‌سازی

زمستان ۹۵

فاز ۲

در این فاز شما مانند فاز قبل تعدادی فایل RSS دریافت کرده و با پردازش آن‌ها را در حافظه خود ذخیره می‌کنید و علاوه بر آن یک پیش‌پردازش روی خبرها نیز انجام می‌دهید.

زمان‌بندی

مهلت انجام این فاز از ۱۶ تا ۳۰ دی است.

ساختار فایل‌ها

شما باید در این فاز تنها یک فایل برای انجام پیش‌پردازش روی خبرهایی که دریافت می‌کنید اضافه کنید. یعنی تمام کدهای این فاز شما در یک فایل خواهد بود و باید این فایل را بعلاوه دو فایلی که از فاز قبلی ساخته‌اید در کوئرا آپلود کنید.

به‌روزرسانی داده‌ها

در ابتدا باید مجموعه‌های url جدیدی که روی سایت قرار داده شده‌است را مانند فاز قبل با فرستادن یک درخواست GET به `fop-project.ir/news/get-urls` دریافت کنید فقط برای اینکه url های فاز ۲ را بگیرید باید در قسمت query-params درخواست خود phase را برابر ۲ قرار دهید. برای روشن شدن موضوع به این نکته دقت کنید: شما در قسمت های قبل در قسمت data درخواست POST خود اطلاعاتی می‌فرستادید مثلاً `answer=salam & member=۱` یا یک رشته json برای مشخص کردن خبر. در درخواست GET قسمت data وجود ندارد و شما باید داده‌های خود را در قسمت query-params وارد کنید که برخلاف درخواست post به انتهای url می‌چسبد. یعنی الان آدرسی که شما باید به آن cURL بزنید برابر است با `www.fop-project.ir/?phase=۲` احتمالاً در هنگام سرچ کردن با گوگل هم دیده‌اید گوگل عبارت شما را در انتهای url جستجو می‌جسباند و از همین روش استفاده می‌کند. درواقع با ظاهر شدن ؟ در انتهای url، قسمت query-params شروع می‌شود. url ها را دریافت و مانند فاز قبل فایل‌های rss مربوط به هر خبر را مجدداً دریافت و ذخیره کنید (دیگر نیازی به فرستادن اطلاعات پردازش شده به سایت نیست. url های خبر در این فاز جدید است و ربطی به فاز قبل ندارد).

پی‌نوشت: احتمالاً الان متوجه شدید چرا درخواست های جستجوی گوگل از نوع get است ولی وقتی مثلاً در یک صفحه شما لاگین می‌کنید درخواست post می‌فرستد و رمز عبور و نام کاربری شما به انتهای url نمی‌چسبد! این ناشی از تفاوت فلسفی get و post است که اولی برای دریافت اطلاعات است و دیگری برای فرستادن اطلاعات و انجام یک کار یا تغییر روی سرور است. برای ساده‌تر شدن کار شما فایل‌های RSS ساده‌تر شده‌اند و هر یک دربردارنده تقریباً ۱۰۰ خبر هستند. به تغییرات فایل‌های RSS دقت کنید. مثلاً اکنون هر خبر تنها یک دسته‌بندی دارد.

دسته‌بندی‌های هدف

اگر کمی در فایل‌های RSS بگردید متوجه این نکته می‌شوید که تنها ۷ دسته‌بندی برای خبرها وجود دارد. هدف ما در این فاز این است که با انجام عملیات‌هایی روی خبرهای دسته‌بندی شده‌مان، بتوانیم به اطلاعاتی برسیم که بتوان با آن‌ها یک خبر که دسته‌بندی آن را نمایانیم گرفته و بگوییم که دسته‌بندی آن چیست.

چیزهایی که می‌توان بیش از همه روی آن‌ها حساب کرد، توضیحات (description) و تیتیر (title) هر خبر هستند که هر کدام شامل تعدادی کلمه هستند که به نوعی نماینده متن خبر محسوب می‌شوند. در

واقع ما قصد داریم تا با انجام پیش‌پردازش روی این بخش‌ها (تیترو توضیحات) به اطلاعات مفیدی راجع به ارتباط کلمه‌ها و دسته‌بندی خبرها برسیم تا بتوان با بررسی کلمه‌ها موجود در تیترو توضیحات یک خبر که دسته‌بندی آن را نمی‌دانیم دسته‌بندی آن را مشخص کرد.

استفاده از داده‌ساختارهای مختلف

برای این که برنامه شما در زمان معقولی اجرا شود، باید از داده‌ساختارهای مختلف مانند درخت پیشوندی (که در سوال اول تمرین ۹ تان نیز باید آن را پیاده‌سازی کنید) در مراحل آتی استفاده کنید (اگر نگاهی به شیوه کار این داده‌ساختار بیاندازید دلیل کاهش زیاد زمان اجرا هنگام استفاده از آن را می‌فهمید). البته باز هم اجباری برای استفاده از این داده‌ساختار ندارید و می‌توانید از روش‌های دیگری مانند درهم‌سازی نیز استفاده کنید اما استفاده از ترای (همان درخت پیشوندی) به شما توصیه می‌شود! نیز فراموش نکنید که سرعت اجرای برنامه شما نیز یکی از پارامترهای نمره‌دهی به پروژه است.

مدل‌سازی و حل مسئله

مسئله ما این است که می‌خواهیم تعدادی خبر که دسته‌بندی آن‌ها را نمی‌دانیم گرفته و با بررسی توضیحات آن‌ها دسته‌بندی‌شان را مشخص کنیم. همان‌طور که قبلاً هم گفتیم برای این کار یک پیش‌پردازش روی متن خبرهایی که دسته‌بندی آن‌ها را می‌دانیم انجام می‌دهیم تا به اطلاعاتی راجع به ارتباط بین کلمه‌ها موجود در متن و دسته‌بندی خبر دست‌یابیم.

(دقت کنید که چون به متن خبرها دسترسی نداریم با تیترو توضیحات آن‌ها کار می‌کنیم.)

در ادامه هر جا به کلمات موجود در خبر اشاره کردیم منظور کلمات تیترو توضیحات خبر در فایل‌های RSS است.)

برای این کار هر خبر را مانند بسته‌ای از کلمه‌ها مدل می‌کنیم! یعنی فقط این اطلاعات را ذخیره می‌کنیم که چه کلمه‌های در تیترو توضیحات این خبر با این دسته‌بندی آمده‌اند و هر کدام چند بار. تا بتوانیم با بررسی تعداد تکرار یک کلمه در خبرهای مربوط به یک دسته‌بندی، تعداد تکرار آن در خبرهای مربوط به بقیه دسته‌بندی‌ها و ... به اطلاعات مفیدی راجع به کلمه‌ها و ارتباطشان با دسته‌بندی‌ها برسیم. البته در این مدل‌سازی اطلاعاتی از خبر مانند توالی کلمه‌ها از دست می‌رود اما در عوض میتوان بسیار ساده‌تر با کلمه‌ها موجود در هر خبر کار کرد.

بررسی ارتباط کلمه‌ها و دسته‌بندی‌ها

کلمه‌های هر خبر (کلمه‌ها موجود در تیترو توضیحات خبر) هر کدام تا حدی مرتبط با موضوع و دسته‌بندی خبر هستند. مثلاً اگر دسته‌بندی یک خبر business باشد، در تیترو توضیحات آن خبر کلمه‌های مثل stocks و bank و sales وجود دارند که ارتباط زیادی با دسته‌بندی دارند و همچنین کلمه‌های مانند new ، for و the وجود دارند که ربط خاصی به دسته‌بندی خبر ندارند. همین‌طور ممکن است نام یک کمپانی گمنام نیز در این کلمه‌ها آمده باشد و چون در جای دیگری نیامده کمکی به ما در دسته‌بندی سایر خبرها نمی‌کند. در مجموع به نظر می‌رسد که کلمه‌های هستند که در تشخیص دسته‌بندی یک خبر از بقیه مفیدتر هستند و می‌توان با ملاک قرار دادن آن‌ها دسته‌بندی یک خبر را حدس زد. یعنی اگر یک خبر که دسته‌بندی آن را نمی‌دانیم به ما بدهند، می‌توانیم با بررسی بودن یا نبودن و تعداد تکرار آن کلمات در آن خبر حدس مناسبی درباره دسته‌بندی آن خبر بزنیم.

اگر بخواهیم علمی‌تر و عملی‌تر (!) صحبت کنیم ما n (که n یک متغیر است) کلمه‌ی موثر در تشخیص دسته‌بندی خبرها را پیدا می‌کنیم (توضیحات بیشتر در بند بعدی داده شده است). سپس هر خبر را با یک بردار به طول n مدل می‌کنیم، برداری که مؤلفه i ام آن نشان دهنده وزن کلمه i ام از n کلمه موثر یادشده

در خبر مورد نظر است. واضح است که دو خبر با موضوع مشابه برداری شبیه به هم دارند. سپس میتوان با مقایسه بردار یک خبر که دسته بندی آن را نمیدانیم با بردار هر دسته بندی (برای هر دسته بندی نیز می توان مشابه یک بردار تعریف کرد) میزان ارتباط آن خبر با دسته بندی را تخمین بزنیم. اگر چیز زیادی از پاراگراف قبلی نفهمیدید، نگران نشوید! چون در این فاز تنها باید آن n کلمه مؤثر در تشخیص دسته بندی را پیدا کنید و بقیه کار مانند به دست آوردن وزن یک کلمه در یک خبر و به دست آوردن بردار هر دسته بندی در فاز بعد مفصل تر شرح داده شده و پیاده سازی می شود.

پیش پردازش

نکته: همان گونه که قبلا نیز اشاره کردیم در ادامه هر جا صحبت از وجود یک کلمه در یک خبر بود منظور همان وجود آن در تیترو توضیحات مربوط به هر خبر در فایل های RSS است. همان طور که گفته شد باید در ابتدا با انجام عملیاتی کلمه ها تیترو توضیحات خبرها را در داده ساختارهایی که طراحی کرده اید ذخیره کنید. برای انجام این کار ابتدا علامت های نگارشی و حروف غیر از حروف الفبای انگلیسی را از تیترو توضیحات هر خبر حذف کرده و در ادامه تمامی بزرگ را به حروف کوچک زبان انگلیسی تبدیل کنید تا بتوان بهتر و ساده تر آن ها را در داده ساختارها ذخیره کرد. حال چگونه باید n کلمه مؤثر در تشخیص دسته بندی ها را پیدا کرد و اصلا این عدد n باید چند باشد؟ در واقع عدد n نیز یکی از پارامترهای مسئله است و می توانید با تغییر دادن آن و بررسی کلمه ها به دست آمده و نتایج قسمت های بعد، آن را تغییر داده و بهتر کنید (فراموش نکنید که در طی مراحل آتی از داده ساختارها که در اول متن فاز نیز به آن اشاره کرده بودیم استفاده کنید...). برای پیدا کردن این کلمه ها باید در ابتدا کل کلمه ها موجود در تیترو توضیحات خبرها را بررسی کنیم. با انجام سه گام زیر روی کلمه ها به n کلمه مؤثر مورد نیاز می رسیم:

- حذف کلمه ها بسیار کم تکرار: اگر یک کلمه در تعداد بسیار کمی از خبرها آمده باشد ملاک معتبری برای مقایسه محسوب نمی شود. مثلا اگر ۲۰۰۰ خبر داشته باشید. کلمه ای که کم تر از ۳ بار تکرار شده باشد اهمیت چندانی ندارد. می توانید با تغییر کران بالای تعداد تکرار (عدد ۳ در این مثال) و بررسی خروجی مقدار مناسب آن را پیدا کنید.
- حذف کلمه ها بسیار پرتکرار: همینطور اگر یک کلمه در تعداد بسیار زیادی از خبرها آمده باشد (مثلا در بیش از سه چهارم خبرها)، نمی تواند ملاک خوبی برای مقایسه باشد. این کار برای حذف کلمه های مانند the و and و for انجام می شود. می توانید با تغییر کران پایین این تعداد و بررسی این که چه کلمه های حذف می شوند مقدار مناسب آن را پیدا کنید.
- انتخاب n کلمه با تاثیر بیشتر در بین کلمه ها باقی مانده (کلماتی که وجود یا عدم وجودشان در یک خبر تاثیر زیادی در حدس ما نسبت به دسته بندی آن خبر دارد).

در نظر داشته باشید که تغییر هر یک از کران ها و متغیرهای بالا تاثیرات مثبت و منفی خاص خودش را دارد و باید با تغییر دادن مقادیر آن ها و بررسی خروجی به مقدار بهینه آن ها برسید. حال باید معیاری برای مقایسه تاثیر یک کلمه در مشخص کردن دسته بندی یک خبر تعیین کنیم. یعنی تابعی طراحی کنیم که با قرار دادن کلمه در آن تابع آن تابع به ما عددی بدهد که نشان دهنده اهمیت آن کلمه در زمینه تعیین دسته بندی باشد (دقت کنید که این تابع باید به طور کلی اهمیت یک کلمه را در تعیین دسته بندی خبر مشخص کند، نه این که فقط اهمیت آن کلمه را نسبت به یک دسته بندی معین مشخص تعیین کند. یعنی شما باید در تابع خود به نحوی نسبت کلمه دریافتی را با تمام دسته بندی ها مشخص کنید، نه فقط یک دسته بندی).

طراحی این تابع به عهده شماست! شما باید با بررسی این که چه عواملی وجود داشتن یا نداشتن یک کلمه در یک خبر را برای تعیین دسته بندی آن مهم می کنند و بررسی آن عوامل در خبرهایی که به شما

داده شده است اهمیت کلمه ها را تعیین کنید و n عدد از مهم ترین آن ها را نیز بیابید. اما برای راهنمایی به فاکتورهای زیر درباره تابع مد نظر که یک کلمه را به عنوان ورودی می گیرد توجه کنید:

- تعداد خبرهای هر دسته بندی.
- تعداد خبرها با دسته بندی C که کلمه ورودی در آن ها هست، به ازای هر دسته بندی C موجود. همین طور تعداد خبرهایی که دسته بندی آن ها C است و کلمه ورودی در آن ها هست نسبت به تعداد کل خبرهایی که کلمه W در آن ها هست، باز هم به ازای هر دسته بندی C موجود.
- تعداد خبرها با دسته بندی C که کلمه ورودی در آن ها نیست، به ازای هر دسته بندی C موجود. همین طور تعداد خبرهایی که دسته بندی آن ها C است و کلمه ورودی در آن ها نیست نسبت به تعداد کل خبرهایی که کلمه W در آن ها نیست، باز هم به ازای هر دسته بندی C موجود.

می توانید این راهنمایی ها را به عنوان سرنخ داشته باشید، از خلاقیت خودتان استفاده کنید و با بررسی نتیجه ای که به ازای استفاده از تابعی که طراحی کرده اید به دست می آید (n کلمه ی نهایی) تابع مناسب را برای انتخاب کلمات موثر بیابید! می توانید این نکته را هم در نظر داشته باشید که چون n کلمه ی موثر برای تعیین دسته بندی خبر های ورودی استفاده می شوند باید در آن ها چند کلمه کلیدی مربوط به هر کدام از دسته بندی ها باشد...