

بسمه تعالی

استخراج واحدهای اندازه‌گیری و کمیت‌ها از متن

گزارش تمرین 1

مرضیه نوری - محمد دل‌خواه

رویکرد کلی

در این پروژه، گروه ما به ترک استخراج مقادیر و کمیت‌ها پرداخت. از آنجا که انواع کلمات و واحدهای این مقوله بسیار زیاد است، تیم ما ابتدا به شناسایی الگوهای بکارگیری متفاوت در جملات برای استخراج کمیت‌ها پرداخت. این الگوها محور اصلی پیاده‌سازی و معماری کد را تشکیل دادند. با این حال با پخته شدن شناخت از این الگوها و پیشرفت پیاده‌سازی، محتوای آن‌ها نیز دستخوش تغییر شد. برخی از الگوها پیش از پیاده‌سازی یا در حین آن با یکدیگر ادغام شدند و برخی دیگر اضافه شدند. همچنین برخی به علت دشواری و دور بودن از بکارگیری مستقیم یک کمیت پیاده‌سازی نشدند.

رویکرد اصلی در یافتن این الگوها چنان بود که به علت تعدد کلمات درگیر در این پروژه، آن‌ها را بر اساس نوع کلی کلمات و عبارات در الگوها تفکیک کنیم و جزییات امکان ترکیب آن‌ها با یکدیگر را به صورت داده‌محور پیاده سازی کنیم. این موضوع در ادامه بیشتر توضیح داده خواهد شد.

بر اساس نوع کلی کلمات و عبارات بکار رفته شده در الگوها، دسته‌های عدد، واحد، صفت، کلیدواژه، ادات سوال و کلمات معین‌گر تعداد شناسایی شد. در بخش توضیح داده‌ها به شرح هر یک خواهیم پرداخت.

D	C	B	A
Done	مثال	الگو	ID
Yes	3 کیلوگرم برنج خریدم یک مائتین به وزن 3 کُن محمد را زیر گرفت باتری خود را 85صم وات شارژ کرد	عدد + واحد + [ایتم]	0
Yes	خانه اتقی به ضلع 4 دارد	کلیدواژه + عدد	1
Merged	مائتین با سرعت زیاد تصالف کرد	کمیت + صفت	2
Yes	رونخاهای با عرض کم دینم	کلیدواژه + صفت	3
Yes	سرع از کترم گشتت	صفت	4
-	یک کیلو شیرینی خرید و نصف آن را خورد		5
Yes	سرعت حرف اول را می‌زند سرعت گرفت	کلیدواژه	6
Merged	نیروی گرانشی اصطکاک عمودبر سطح کششای الکترومغناطیسی مقاومت هوا	کمیت + نوع کمیت	8
Merged	بیروی اصطکاک زیاد	کمیت + نوع کمیت + صفت	9
-	مثال برقی از کترم گشتت	تشبیه بدون وجه شبه صریح	10
-	وزنهای به سذگینی یک فیل را بلند کرد.	به + کلیدواژه + ی + اسم	11
-	مثال فیل سدگین بود	مثال + اسم + کلیدواژه	12
Merged	چند کیلو داریم؟	ادات سوال + واحد	13
Yes	چند صد ژول انرژی	ادات سوال + دهیمی + واحد + [ایتم]	14
			15

معماری و ساختار کد

همان طور که در فایل ipynb پروژه قابل مشاهده است، کد از 2 نوع قسمت تشکیل شده:

- MeasurementExtractor کلاس
- توابع تشخیص الگو

البته این تفکیک تنها به علت تست و کارکردن راحت‌تر با کد و استفاده از امکانات سول‌های jupyter است و به راحتی می‌توان تمام توابع تشخیص الگوها را به داخل کلاس MeasurementExtractor منتقل نمود.

مطابق خواست دستورالعمل تمرین، تابع اصلی کلاس run بوده که با گرفتن یک رشته در خروجی قطعه‌هایی که به کمیت‌ها اشاره دارند را به همراه اطلاعات قابل استخراج آن‌ها به صورت آرایه‌ای از jsonها باز می‌گرداند. کار اصلی تابع run اجرا و ادغام توابع تشخیص الگو و در نهایت بازگرداندن خروجی نهایی است.

کار اصلی دیگر کلاس MeasurementExtractor خواندن دادگان مورد نیاز از فایل‌ها و آماده‌سازی آن‌ها برای استفاده است. همچنین توابع کمکی مورد استفاده در برخی الگوها در آن قرار دارد.

داده‌ها

همان طور که ذکر شد، در این پروژه تلاش شده رویکرد تا جای ممکن data-driven باشد و داده‌ها از کد تفکیک شوند. تشخیص عدد و گروه اسمی بر عهده کتابخانه‌های خارجی `parsi.io` و هضم قرار داده شده و دیگر المان‌های محوری موجود در الگوها، یعنی واحد، صفت، کلیدواژه، ادات سوال و کلمات معین‌گر تعداد، همه به صورت فایل‌های CSV مجزا در کنار پروژه قرار داده می‌شوند. تا بدون نیاز به شناخت کد بتوان آن‌ها را تکمیل کرد و ارتقا داد.

واحدها

واحدها از سایت باحساب استخراج شده و در یک فایل CSV قرار داده شده که هر خط آن متعلق به یک نوع کمیت، و عناصر هر خط به صورت جفت‌های «کلمه واحد» و «نماد واحد» هستند.

	A	B	C	D	E	F	G	H
1	length	m متر		cm سانتیمتر		inch اینچ		mm میلیمتر
2	mass	Ton تن		kg کیلوگرم		lb پوند		g گرم
3	pressure	pa پاسکال		bar بار		psi پی اس آی		kpa کیلو پاسکال
4	volume	L لیتر		mL میلی لیتر		متر مکعب		سانتیمتر مکعب
5	temperature	C سانتیگراد		F فارنهایت		K کلوین		R رانکین
6	area	متر مربع		سانتیمتر مربع		میلی متر مربع		کیلومتر مربع
7	speed	m/s متر بر ثانیه		fps فوت بر ثانیه		km/h کیلومتر بر ساعت		mph مایل بر ساعت
8	force	N نیوتن		gf گرم-نیرو		kgf کیلوگرم-نیرو		lbf پوند-نیرو
9	energy	J ژول		mJ میلی ژول		μJ میکرو ژول		kJ کیلو ژول
10	power	W وات		mW میلی وات		kW کیلو وات		MW مگا وات
11	torque	N.m نیوتن-متر		ft.lbf فوت-پوند نیرو		m.Kgf متر-کیلوگرم نیرو		m.Kf متر-کیلو پوند
12	time	s ثانیه		ks کیلو ثانیه		ms میلی ثانیه		μs میکرو ثانیه
13	density	t/m³ تن بر متر مکعب		kg/m³ کیلوگرم بر متر مکعب		kg/L کیلوگرم بر لیتر		kg/dr کیلوگرم بر دسی متر م
14	frequency	Hz هرتز		rps دور بر ثانیه		rpm دور بر دقیقه		rph دور بر ساعت
15	degree	Rad رادیان		Deg درجه		Grad گراد		arc s ثانیه قوسی
16	acceleration	متر بر مجذور ثانیه		کیلو متر بر مجذور ثانیه		سانتی متر بر مجذور ثانیه		فوت بر مجذور ثانیه
17	debi	تن بر ثانیه		تن بر دقیقه		تن بر ساعت		تن بر روز
18	debi-v	متر مکعب بر ثانیه		متر مکعب بر دقیقه		متر مکعب بر ساعت		متر مکعب بر روز
19	data-storage	Bit بیت		kb کیلو بیت		Kib کیبی بیت		Mb مگا بیت
20	data-transfer	Bit/s بیت بر ثانیه		kb/s کیلو بیت بر ثانیه		Kib/s کیبی بیت بر ثانیه		Mb/s مگا بیت بر ثانیه
21	ratio	درصد	%					
22	wildcard	واحد						

صفات

صفات از رده «صفات فارسی» سایت `wiktionary` استخراج شده‌اند و در میان آن‌ها صفاتی که مرتبط با کمیت‌ها و قابل استفاده در الگوها هستند در ستون Useful با عدد 1 به عنوان پرچم مشخص گشته‌اند. همچنین در میان صفات مرتبط، همه در تمام الگوها قابل استفاده نیستند برای همین این که در چه الگویی امکان استفاده از آن‌ها وجود دارد با یک پرچم در ستون آن الگو مشخص شده است. در نهایت برای الگوی 4، از آنجا که نوع کمیت باید از صفت استنتاج شود، در یک ستون نوع کمیت ارجح برای آن صفت مشخص شده است. شایان ذکر است همه صفات در `wiktionary` موجود نبودند و برخی مانند خنک، گرم و ... به صورت دستی توسط تیم به آن اضافه گشته‌اند.

	A	B	C	D	E
1	adjective	useful	valid for pattern 4	pattern 4 type	valid for pattern 3
2	آبی	1			
3	آرام	1		1 speed	1
4	آستان	1		1 length	
5	آستین دراز	1			
6	آشکار	1			
7	آکبند	1		1 time	
8	آهسته	1		1 speed	1
9	ابتدایی	1		1 time	1
10	ارزان	1			
11	افسارگسیخته	1			1
12	البوه	1			
13	اندک	1			1
14	انگیزش	1		1 energy	
15	اونجا	1			
16	ایستاده	1		1 speed	

کلیدواژه‌ها

کلیدواژه‌ها نیز مشابه صفات بوده با این تفاوت که همه آن‌ها به صورت دستی توسط اعضای تیم نوشته شده‌اند. در کلیدواژه‌ها نیز مشابه صفات الگویی نیازمند تفکیک کلمات موجود به قابل استفاده و غیر قابل استفاده بود که با روش پرچم گذاری انجام شد.

acceleration	افزایش سرعت	1	
acceleration	کاهش سرعت	1	
acceleration	تخیز سرعت	1	
acceleration	میزا	0	
acceleration	ترمز	1	
data-storage	حجم	1	
data-storage	اندازه	1	
data-storage	بزرگ	0	
data-storage	کوچک	0	
data-storage	حفظه	1	
mass	وزن	1	
mass	سنگینی	1	وزنمای به سنگینی یک فل را بلند کرد.
mass	سنگین	0	بسته سنگین بود.
mass	سبکی	1	چیزی به سبکی زیر هوا معلق بود.
mass	سبک	0	بسته سبک بود.
mass	ترازو	1	ترازو روی عدد 2 ایستاد.
mass	جرم	1	جسمی با جرم بالا سقوط کرد. جسمی به جرم دو کیلوگرم سقوط کرد.
temperature	دما	1	
temperature	نمای	1	
temperature	داغ	0	آب داغ بود
temperature	سرد	0	آب سرد بود.

کلمات معین‌گر تعداد

در این فایل علاوه بر پیشوندهای یونانی، کلمات خاص مانند جین و دوجین و کلمات اعداد دودویی (کیبی، گیبی و ...) قرار داده شده است. بنا بود این کلمات از ابتدای واحدها جدا شده و به صورت مجزا تشخیص داده شوند تا با توجه به عدد متناظر آنها بتوان در تبدیل واحد از آنها استفاده نمود. متأسفانه با توجه به اتمام وقت این کار پیاده‌سازی نشد.

میلیدانوم		0.000000001
نریلیونوم		0.000000000001
کواندالیونوم		0.0000000000000001
کوئینتالیونوم		0.0000000000000000001
سکستالیونوم		0.0000000000000000000001
سپتالیونوم		0.000000000000000000000001
جفت	Pair	2
کیبی	Ki	1024
می	Mi	1048576
گیبی	Gi	1073741824
تی	Ti	1099511627776
پی	Pi	1125899906842624
اگزیبی	Ei	1152921504606846976
زی	Zi	1180591620717411303424
یوی	Yi	1208925819614629174706176
جین		6
دوجین	dozen	12

فایل‌های کمکی

در این فایل‌ها بعضی کلمات بجای نوشته شدن در کد قرار گرفته اند مانند ترجمه نوع کمیت‌ها و ادات سوال.

	A	B
1	english	persian
2	length	طول
3	mass	وزن
4	pressure	فشار
5	volume	حجم
6	temperature	دما
7	area	مساحت
8	speed	سرعت
9	force	نیرو
10	energy	انرژی
11	power	توان
12	torque	گشتاور
13	time	زمان
14	density	چگالی
15	frequency	فرکانس
16	degree	زاویه
17	acceleration	شتاب
18	debi	شارش جرمی
19	debi-v	شارش حجمی
20	data-storage	ذخیره دیجیتال
21	data-transfer	انتقال داده
22	wildcard	علم

	A
1	چند
2	چقدر
3	چند
4	اینقدر
5	اینقدر
6	آنقدر
7	آنقدر
8	همینقدر
9	همینقدر
10	چه مقدار
11	چندین

شناسایی آیتم

با توجه به آن که نقش کلمات در جملات بسیار متفاوت بوده و گاهی میان آنها حرف اضافه یا علایم نگارشی خاصی ظاهر نمی‌گردد، تشخیص مسقیم آیتم با استفاده از عبارات منظم همیشه ممکن نیست. برای این کار

گروه ما از قابلیت chunker کتابخانه هضم استفاده نمود تا گروه‌های اسمی را تفکیک کند و هنگامی که در ادامه یک کمیت اسمی آمده باشد، تنها در صورتی به عنوان آیتم شناسایی شود که جزو همان گروه اسمی باشد. با این حال متاسفانه مدل chunker کتابخانه هضم همیشه درست کار نمی‌کند و باعث بروز خطا می‌شود. استفاده درست از علایم نگارشی مانند نقطه در انتهای جمله، خطای این کتابخانه را کاهش می‌دهد اما به صفر نمی‌رساند.

عدم وابستگی به کد و پیاده‌سازی

از مزیت‌های نحوه پیاده‌سازی این پروژه توسط گروه ما، تفکیک مناسب داده و رفتار است به طوری که در صورت نیاز به افزودن کمیت‌های کاملاً جدید، هیچ نیازی به تغییر کد نبوده و کافی است اطلاعات لازم در فایل‌های csv اضافه شود. با این اتفاق تمام الگوهای موجود کمیت جدید را شناسایی کرده اطلاعات لازم را از نمونه‌های آن در متن داده شده استخراج می‌کنند. تنها در صورتی تغییر کد نیاز است که الگوی جدیدی مد نظر باشد که برای این کار کافی است تابعی برای آن نوشته شده و در run کلاس MeasurementExtractor صدا زده شود.