

## گزارش پروژه ناشناس کردن اطلاعات شخصی

نام خانوادگی اعضای گروه: نفیسه نیک اقبال - محمد هادی حاجی حسینی - الیاس اسماعیلی

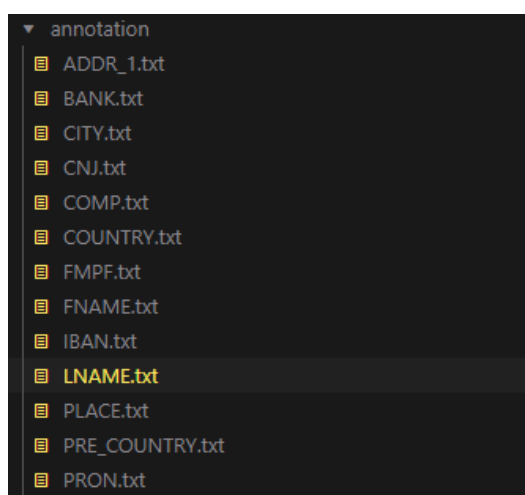
پوشه های اصلی پروژه :

پوشه annotation :

داخل این پوشه دیتاست های ما مانند لیست اسامی، لیست نام های خانوادگی، لیست نام شرکت ها، لیست اسامی شهرها و کشورها وجود دارد.

علاوه بر موارد ذکر شده واحدهای سازنده هر Rule ها مانند کلمات خاصی که در آدرس ها استفاده می شوند مانند کلمات شهر، خیابان، جنب، روبه رو، علائم ربط مثل کاما و ویرگول، ضمائر و... آورده شده اند که ما از فایل های داخل این پوشه استفاده می کنیم تا Pattern ها را بسازیم که در حقیقت بخشی از واحدهای سازنده pattern ها این annotation ها می باشند.

لیست فایل های داخل پوشه annotation به شکل زیر است:



فایل ADDR\_1 : کلمات کلیدی موجود در آدرس ها مانند خیابان، کوچه، جنب، بزرگراه و... می باشد.

فایل BANK : شامل 6 شماره معتبر اول شماره های کارت می باشد.

فایل CITY : شامل اسامی تمامی شهرهای ایران و دنیا می باشد.

فایل COUNTRY : شامل اسامی تمام کشورها می باشد.

فایل PLACE : مکان های مختلف مثل سینما، قوه قضائیه، رستوران و ... می باشد.

فایل IBAN : شامل regex های معتبر کشورهای مختلف برای IBAN می باشد.

فایل COMP : نام کلیه شرکت هاهم انگلیسی و هم فارسی در این فایل موجود می باشد. که با crawl سایت دانشکار و jobinja بدست آمده اند.

فایل FNAME و LNAME : یک مجموعه از اسامی فارسی و فامیل های فارسی رایج می باشد که بخشی از آن ها با crawl سایت ویکی پدیا بدست آمده است.

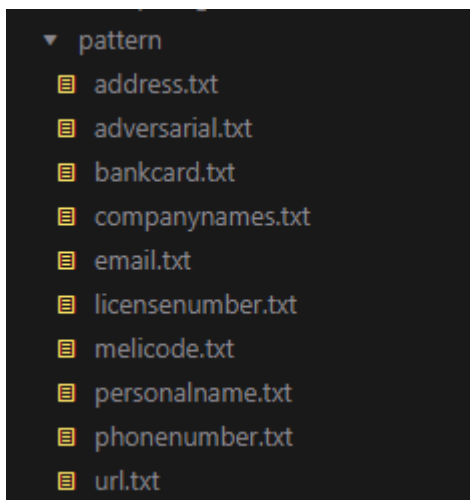
فایل PRON : شامل تمام ضمائر جدا و چسبیده فارسی می باشد.

فایل CONJ : شامل ویرگول و کاما و - می باشد.

فایل FMPF : این فایل تمامی پسوندهای معتبر برای فامیل مانند فر، نژادو... می باشد

## : Pattern

ما داخل این پوشه برای هر category از اطلاعات که می خواهیم اطلاعات آن را مخفی کنیم یک فایل در نظر می گیریم و در داخل pattern ها ما با استفاده از annotation ها قانون های خود را می نویسیم که در نهایت این rule ها تبدیل به regex می شوند.



به عنوان مثال فایل personalname شامل قانون هایی است که برای شناسایی نام افراد می خواهیم داشته باشیم. که هر خط از داخل این فایل به یک rule تبدیل می شود و داخل فایل هایی که در ادامه توضیح می دهیم به regex تبدیل می شود.

در خط اول و دوم می گوید هر جا یکی از اسم ها یا فامیل های داخل فایل های FNAME یا LNAME آمد را در نظر بگیر و خط بعدی نیز می گوید بین اسم و فامیل می تواند علائم ربط مثل - نیز بیاید و در خطوط بعدی می گوییم اگر کلمات اسم، نام، نام خانوگی و فامیل آمدند و بعد آن ها یک ضمیر منفصل یا متصل آمد کلمه بعدی آن ها را نیز به عنوان اسم در نظر بگیرد و این باعث می شود که اگر جمله ای مانند "اسم من جیسون اکبری است" را کلمه جیسون را به عنوان اسم برگرداند و هم چنین FNAME<>FMPF برای ما از روی اسامی و پسوندهای تعریف شده فامیل می سازد و علامت <> یعنی صفر فاصله یا بیشتر می توانند داشته باشند. به

عنوان مثال نام محمد را داریم و از روی آن فامیل‌های محمدی، محمدی فر، محمدپور و.... را می‌سازد. که قانون‌های توضیح داده شده در شکل زیر آورده شده اند:

```
# Personal Names
FNAME
LNAME
FNAME<>CNJ<>LNAME
اسم PSPACEPRON ONE_W
نام PSPACEPRON ONE_W
فامیل PSPACEPRON ONE_W
نام خانوداگی PSPACEPRON ONE_W
نام خانوداگی PSPACEPRON ONE_W
FNAME<>FMPF
FNAME<>FMPF<>FMPF
```

فایل pattern\_to\_regex :

در این فایل ما دو کلاس داریم :

1. کلاس Annotation :

تمام annotation ها را از روی فایل‌ها خوانده و یک دیکشنری ساخته که کلید آن برابر اسم فایل می‌باشد و value آن OR تمام اطلاعات خوانده شده از فایل‌ها می‌باشد و هم چنین اگر regex ای هم بخواهیم اضافه کنیم ک داخل متن فایل‌های annotation نبوده و بخواهیم در داخل pattern ها از آن استفاده کنیم را مستقیم اینجا تعریف کرده و به دیکشنری خود اضافه می‌کنیم.

2. کلاس Pattern :

در این کلاس ما فایل‌های پوشه Pattern را خوانده و یک دیکشنری باید درست کنیم که کلید آن شامل نام فایل‌های پوشه Pattern است و value آن شامل تمام regex های قانون‌هایی است که در داخل این فایل‌ها نوشتیم و در حقیقت در داخل این کلاس با استفاده از دیکشنری ساخته شده در کلاس annotation یک دیکشنری جدید می‌سازیم که شامل تمام regex ها ما هست.

فایل Spans :

در این فایل کاری که انجام می‌دهد این است که با توجه به دیکشنری ساخته شده در مرحله قبل span تمام regex های نوشته شده را بدست می‌آورد و اگر دوتا pattern پیدا کرده باشد و بین آن‌ها فقط space باشد

آنگاه span های آن ها را با هم merge می کند. به عنوان مثال در جمله "علی رضا را در مدرسه دیدم" علی را داخل نام ها پیدا می کند رضا را هم پیدا می کند و بعد چون بین آن ها فقط space است آن ها را به هم می چسباند. هم چنین چون اسم آدرس ها گاهی به عنوان اسم اشخاص در نظر گرفته می شود ما در این کلاس اگر یک اسم به عنوان آدرس تشخیص داده شده باشد دیگر به عنوان نام اشخاص در نظر نمی گیریم.

فایل anonymizer :

در این فایل یک کلاس Model تعریف می کنیم و در آن از کلاس pattern استفاده کرده و regex ها را ساخته و بعد از پکیج parstdex نیز استفاده می کنیم تا زمان ها و تاریخ را استخراج کنیم و بعد در نهایت تمام span ها در قالب یک دیکشنری برمی گردانیم و در کلاس example از مدل ساخته شده استفاده می کنیم.

ورودی :

```
sent =
.علی احمدی در شهرستان اصفهان شهر فولادشهر به دنیا آمد.
.علی رضا و زهرا در کشور زیمبابوه زندگی می کنند
.به نظر سخت است france زندگی در کشور
.کشور به منطقه ای گفته می شود که مرز آن با سیاست تعیین شده است
.من در ۱۶ بهمن ۱۳۷۵ به دنیا آمدم
.در تاریخ ۱۶ فروردین به مکه رفتم و از شرکت داده پردازش نیز دیدن کردم
.کد ملی من ۱۱۳۰۳۹۶۷۸۹ است
.کد ملی من است ۱۱۳۰۳۹۶۷۸۹
. شماره کارت اعتباری من 6104337958646987 است که از آن به حساب شما پول می ریزم. من سایت
. است IR069600000001032420000011 شماره شبای من
. است KW81CBKU0000000000001234560101 من IBAN شماره
. ساعت ۸ صبح پرواز به مقصد روستای احمدآباد را دارم
. من در اصفهان، کوشک، خیابان احمدی، کوچه شهید علی علیزاده پلاک ۱۴۳ زندگی میکنم
. این نشانی خیابان بهشتی است
. شماره تلفن من ۰۹۱۲۳۴۵۶۷۸۹ و شماره خانه علیرضا ۰۲۱۳۳۴۵۵۵۶۶ است
. خیابان اهواز شهر اهواز بسیار تمیز است
```

خروجی:

```
نام شخص#> در <#آدرس#> به دنیا آمد#>
نام شخص#> و <#نام شخص#> در <#آدرس#> زندگی میکنند#>
زندگی در <#آدرس#> به نظر سخت است
کشور به منطقه ای گفته میشود که مرز آن با سیاست تعیین شده است
من در <#تاریخ و ساعت#> <#تاریخ#> به دنیا آمدم
در <#تاریخ و ساعت#> <#تاریخ#> به <#آدرس#> رفتم و از <#اسم شرکت#> نیز دیدن کردم
کد ملی من <#کد ملی#> است
کد ملی#> کد ملی من است#>
سایت <#آدرس سایت#> رو مشاهده کردم و ایمیل کارمند <#اسم شرکت#> <#ایمیل#> را برداشتم
شماره <#تاریخ و ساعت#> <#تاریخ#> من <#اطلاعات حساب#> است
من <#اطلاعات حساب#> است IBAN شماره
تاریخ و ساعت#> <#زمان#> پرواز به مقصد <#آدرس#> را دارم#>
من در <#آدرس#> زندگی میکنم
این نشانی <#آدرس#> است
شماره تلفن من <#شماره تماس#> و شماره خانه <#نام شخص#> <#شماره تماس#> است
آدرس#> بسیار تمیز است#>
```

ما از پایتون 3 استفاده کردیم و از کتابخانه re برای شناسایی regexها و از parstdex برای شناسایی زمان و تاریخ استفاده نمودیم.