

به نام خدا

محمد مهدی آقاجانی

۹۳۳۱۰۵۶

تمرین دوم

استاد : دکتر ممتازی

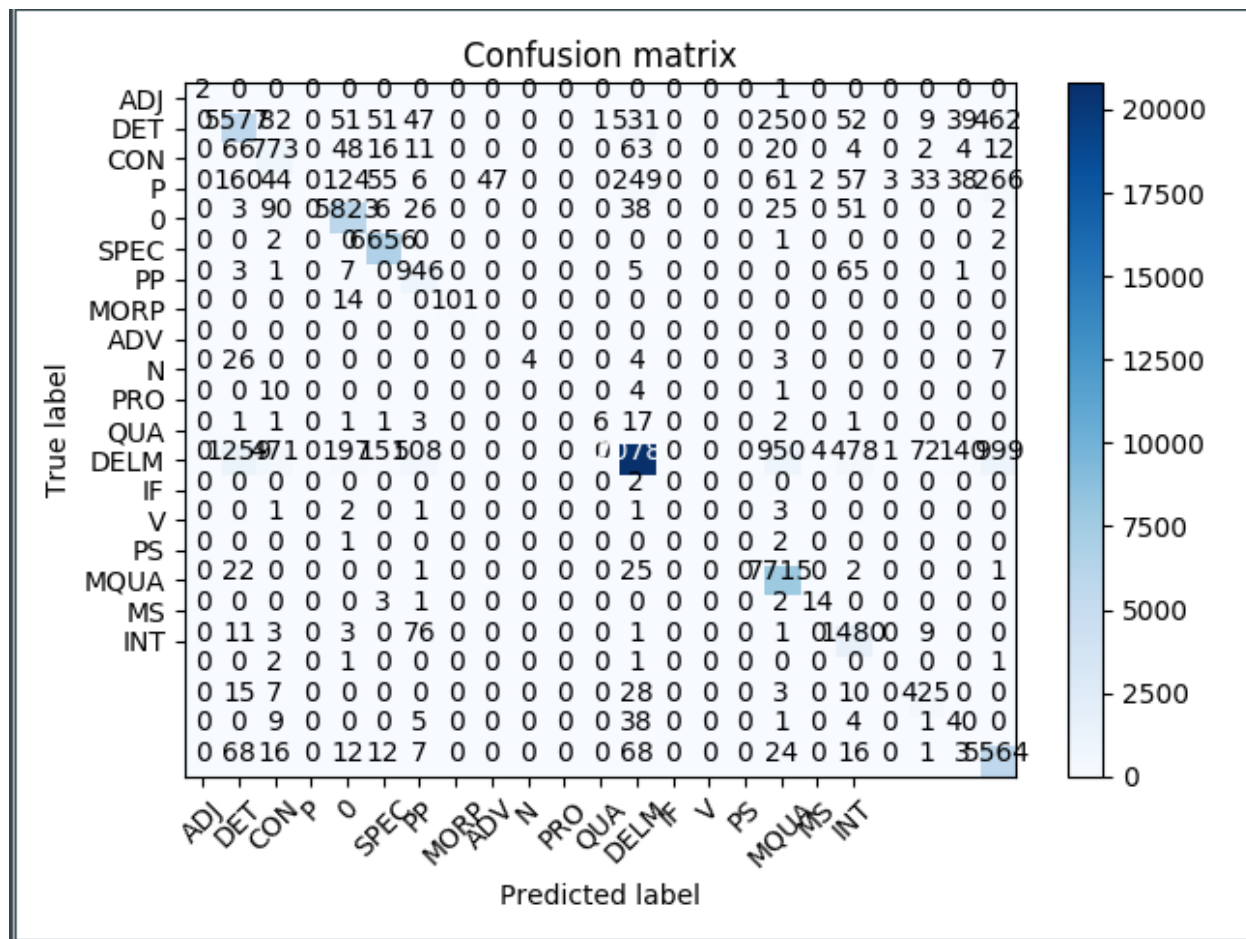
تمرین اول

در این تمرین از ابزار NLTK استفاده شد. برای عمل آموزش از متد `HiddenMarkovModelTagger` استفاده کردیم که آن را با داده هایی که در اختیار داشتیم `train` کردیم و سپس با استفاده از متد `tag` اقدام به تگ گذاری داده ها نمودیم.

نکته بسیار مهم این است که در این نوع یادگیری باید حتما داده ها را بر اساس جمله و نه کلمه له کلمه در اختیار ابزار بگذاریم تا دقت بهتری حاصل شود.

با استفاده از روش بالا دقت ۸۵ درصد به دست آمد که دقت مناسبی میتواند باشد البته می توان با تنظیم پارامترها به دقت های بهتری دست پیدا کرد. ماتریس درهم^۱ این سوال نیز در گزارش آمده است

^۱ Confusion matrix



تمرین دوم

در این تمرین از ابزار Stanford استفاده نمودیم. ابتدا لازم بود که مدل مربوط به خودمان را با داده های آماده شده بسازیم که برای این کار باید یک فایل prop. آماده میشد که مسیر داده و پارامترهای مورد نظر در آن ست شده بود (این فایل درون پوشه بارگذاری شده موجود می باشد). سپس با استفاده از دستور زیر مدل خود را ساختیم :

```
java -cp stanford-ner.jar edu.stanford.nlp.ie.crf.CRFClassifier -prop filename.prop
```

با این دستور یک فایل ser.gz. تولید میشود که باید در کد از آن استفاده کنیم.

سپس در کد با استفاده از متد tag آن ها را دسته بندی کردیم. در این حالت دقت کار ۹۵ درصد شد که قابل قبول است اما نکته مهم توجه به confusion matrix است که در گزارش هم آمده است و در این نمودار میبینیم که خیلی از کلاس ها به کلاس PERSON نگاشت شده اند. البته تا حدودی نشان از این بحث چالش برانگیز می باشد که خیلی از این دسته بندی ها با دسته بندی person اشتراک زیادی دارند زیرا اسامی اشخاص میتواند عمومیت بیشتری نسبت به ما بقی کلاس ها داشته باشد.

