

برای انجام بخش اول این فاز ۴ تابع تعریف کردم تا متن ورودی را به ترتیب نرمال کرده ( مواردی مانند درست کردن تاریخ ، درست کردن فاصله و نیم فاصله ، تبدیل اعداد فارسی به انگلیسی ، تبدیل نسخه فینگلیش به فارسی و ... ) در مرحله بعدی توکن ها را در یک لیست جمع آوری میکند . در مرحله بعدی حروف یا کلمات پرتکرار مثل با - به - از و ... و علائم نگارشی مثل نقطه ویرگول و... را از لیست توکن ها حذف میکند . مرحله آخر کلمات با ریشه یکسان را به یک کلمه تبدیل میکند مثلاً صورت مختلف صرف افعال .

برای دو بخش اول از parsivar و بخش آخر از تابع های hazm استفاده شده است. برای حذف stopword ها از یک فایل تکست استفاده شده که عبارات موجود در آن را با توکن ها مقایسه میکنیم و اگر در لیستمان موجود بود حذف میکنیم.

نتیجه برای متن ورودی "من به تاریخ ۱۱ شهریور به دانشگاه می روم. man be tarikhe 11 shahrivar be danshgah miravam" به شکل زیر است :

```
miravam شهریور به دانشگاه 11 tarikhe من به تاریخ 11 شهریور به دانشگاه می روم . منع به : Normal shode
'شهریور', 'به', 'دانشگاه', '11', 'tarikhe', 'روم', 'من', 'به', 'تاریخ', '11', 'شهریور', 'به', 'دانشگاه', 'می' : Token haye be dast amade
'miravam', 'شهریور', 'دانشگاه', '11', 'tarikhe', 'روم', 'من', 'تاریخ', '11', 'شهریور', 'دانشگاه', 'می' : Hazf kalamate por tekrar
'miravam', 'شهریور', 'دانشگاه', '11', 'tarikhe', 'روم', 'من', 'تاریخ', '11', 'شهریور', 'دانشگاه', 'می' : Hazf kalamate ba rishe yeksan
```

می بینیم که بصورت کلی به جواب خوبی میرسیم مثلاً نیم فاصله در می روم را درست کرده - عدد فارسی را انگلیسی کرده - فینگلیش را تا حد خوبی درست برگردانده اما جاهایی که اعراب دارد مثل اعراب کسره در تاریخ را تشخیص نداده یا man را به جای من ، منع ترجمه کرده .

در بخش تهیه توکن ها هم مشکلاتی مثل مورد تکراری یا اشتباهی در نظر گرفتن کلمه مرکب مثلاً در اینجا به من را یک کلمه حساب کرده (u200c\ برای فاصله بین کلمات مرکب است که در اصل همان نیم فاصله است)

در بخش حذف stopword درست عمل کرده و "به" را حذف کرده در بخش ریشه یابی هم رفتن را به رو تبدیل کرده .

برای مشکلات ذکر شده در بخش های بعدی بصورت دستی حل میشود.