

# تمرین 1

7 آبان

مهدیه ساجدی پور

## 1. مقدمه

هدف این تمرین بررسی مراحل مختلف در پیش پردازش متن فارسی و همچنین آشنایی با ابزارهای مختلف پردازش زبان فارسی هست. سه ابزار مختلف بر روی متن داده شده استفاده می‌شوند تا تفاوت‌های آنها در نحوه ی پیش پردازش مقایسه شود.

## 2. توضیحات

پیش پردازش متون به سه بخش نرمال سازی، ریشه‌یابی و بن‌یابی، و توکن‌بندی تقسیم می‌شود. ما در این تمرین، این مراحل را با ابزارهای هضم، دادما، و پارس‌یوار انجام می‌دهیم. متن آزمایشی به این صورت است:

"

عنوان مقاله: صفحه اصلی

<p/>متن ویکی‌پدیا فارسی

من john هستم و در ب.م.م گیری تخصص دارم!

کلمات عربی مانند اصلاح کاف و یا ی برای توکنایزر ما اهمیت دارند.

ما می دانیم که در تاریخ ۲۰ سپتامبر ۲۰۰۴ (۲۹ شهریور، ۱۳۸۳) مقاله های "ویکی پدیا" در

۱۰۵ زبان به یک میلیون رسید.

که این مقالات شامل زمان های پیشین نیستند.

در ویکی‌پدیا فارسی ممکن است ( گاهی ) فاصله پرانتز ها رعایت نشده باشد، یا حتی ممکن

است درباره محبوب ترین های فارسی صحبت شده باشد.

در اینجا یک ایمیل آزمایشی از من sh@sbu.ac.ir قرار دارد.

برای اطلاعات بیشتر می‌توانید به وبسایت ویکی‌پدیا فارسی به آدرس <http://wikipedia.com>

سر بزنید.

(داخل پرانتز بگویم، این یک متن تستی است. حداقل به من اینطور گفته شده است.)

مجله تایم در گزارش سال ۲۰۰۶ خود، جیمی ویلز را در گروه ۱۰۰ فرد تأثیرگذار سال اعلام کرد.

همچنین در همین سال ویکی پدیای روسی برنده جایزه رانیت (Премия Рунета (روسی: در بخش «دانش و آموزش» شد. این جایزه از طرف دولت اعطا می شود.

همچنین ویکی پدیا جایزه یک میلیون دلاری مدیریت پروژه را از همایش صفاجو دریافت کرد.

پلتفورم اهداف توسعه پایدار United Nations:

چندین پروژه متن-آزاد دارد که وظایف غیردانشنامه ای را انجام می دهند

"

## a. نرمال سازی

نرمال سازی از مراحل ابتدایی در پردازش هر زبانی است. زبان فارسی دارای متونی سخت برای نرمال سازی است که دلایل مختلفی دارد. از جمله، فاصله‌ها ممکن است به طرق مختلفی نمایان شوند؛ به طور مثال کلمه درخت‌ها می‌تواند به چند صورت درختها، درخت‌ها و درخت‌ها نوشته شود. قانون ه کسره در بسیاری از متون رعایت نمی‌شود و ... در نرمال‌سازی فاصله‌های بین کلمات، علائم نگارشی، تکرار در حروف، اختصارها و ... اصلاح می‌شوند تا به یک استاندارد مشترک بین پیش‌پردازش و پردازش که همان یونیکد است برسیم.

## i. متن نرمال شده توسط هضم:

'عنوان مقاله: صفحه اصلی<n\n<p\ متن ویکی\200c\پدیا فارسی> / >n\n\p>n\ من john هستم و در ب. م. م گیری تخصص دارم! \n\ کلمات عربی مانند اصلاح کاف و یا\200c\ی برای توکنایزر ما اهمیت دارند. \n\ما می\200c\دانیم که در تاریخ ۲۰ سپتامبر ۲۰۰۴ (۲۹ شهریور، ۱۳۸۳) مقاله\200c\های «ویکی پدیا» در ۱۰۵ زبان به یک\200c\میلیون رسید. \n\که این مقالات شامل زمان\200c\های پیشین نیستند. \n\در ویکی\200c\پدیا فارسی ممکن است (گاهی) فاصله پراتتز\200c\ها رعایت نشده باشد، یا حتی ممکن است درباره محبوب\200c\ترین\200c\های فارسی صحبت شده باشد. \n\در اینجا یک ایمیل آزمایشی از من sh@sbu. ac. ir قرار دارد. \n\برای اطلاعات بیشتر می\200c\توانید به وبسایت ویکی\200c\پدیا فارسی به آدرس http: // wikipedia. com سر بزنید. \n\ (داخل پراتتز بگویم، این یک متن تستی است. حداقل به من اینطور گفته\200c\شده است.) \n\مجله تایم در گزارش سال ۲۰۰۶ خود، جیمی ویلز را در گروه ۱۰۰ فرد تأثیرگذار سال اعلام کرد. \n\همچنین در همین سال ویکی پدیای روسی برنده جایزه رانت (Премия Рунета) در بخش «دانش و آموزش» شد. این جایزه از طرف دولت اعطا می\200c\شود. \n\همچنین ویکی پدیا جایزه یک\200c\میلیون دلاری مدیریت پروژه را از همایش صفاجو دریافت کرد. \n\پلتفورم اهداف توسعه پایدار United Nations: \n\چندین پروژه متن-آزاد دارد که وظایف غیردانشنامه\200c\ای را انجام می\200c\دهند'



نرمالایزر پارسیوار <p> را به سه کلمه تبدیل میکند اما دیگر ابزار ها آن را یک کلمه در نظر میگیرند.

هیچ کدام کلمات مختصر شده را تشخیص نمی‌دهند.

تنها پارسیوار به قبل علائم نگارشی فاصله می‌افزاید.

دادما از دیگر ابزارها در تشخیص نیم‌فاصله‌ها و تصحیح‌شان ضعیف‌تر عمل می‌کند.

هر سه در تبدیل ی به ی درست عمل می‌کنند اما دادما اعراب را حذف نمی‌کند.

دادما در مرحله نرمالایز کردن قابلیت تبدیل اعداد و url و ادرس ایمیل را به token های مشخص شده دارد.

پارسیوار میتواند تاریخ شمسی را نرمالایز کند.

دادما بر خلاف دو ابزار دیگر قابلیت حذف کاراکتر — در زمان را ندارد.

در پارسیوار فاصله ی پرانتز ها تصحیح نمی‌شود(استاندارد این کتابخانه با بقیه متفاوت است).

تنها دادما میتواند ایمیل و url را تشخیص دهد، بقیه قبل و بعد از . فاصله می‌گذارند.

دادما افعال چندبخشی را تشخیص نمی‌دهد، هضم تنها افعال را به صورت دو بخشی تشخیص می‌دهد، اما پارسیوار میتواند افعال سه بخشی را نیز تشخیص دهد.

هضم برخلاف دو ابزار دیگر، ا را به تبدیل نمی‌کند.

تنها هضم ء را حذف می‌کند.

پارسیوار ترکیب متن-آزاد را به سه توکن تبدیل می‌کند.

به نظر می‌رسد پارسیوار متن نرمال شده‌ی بهتری را تولید می‌کند.

## b. توکن‌بندی

ساده‌ترین راه برای توکن‌بندی استفاده از فاصله به عنوان جدا کننده است. توکن‌بندی به ما کمک می‌کند بتوانیم از ابزارهای ریشه‌یابی و بنیابی استفاده کنیم و نیز روابط بین کلمات در هر جمله را به دست بیاوریم. همچنین همچنین هر توکن به عنوان یک لغت در لغت‌نامه مدل ذخیره می‌شود. بهتر در توکن‌بندی از ابزاری استفاده شود که کلمات چندبخشی را تشخیص داده و به عنوان یک توکن در نظر بگیرد.

## i. توکن‌بندی با هضم

```
[ 'عنوان', 'مقاله', ':', 'صفحه', 'اصلی', '<p>متن',  
'ویکی\200c\پدیا', 'فارسی', '>', '/', '<p', 'من', 'john', 'هستم',  
'و', 'در', 'پ', '.', '4', 'NUM', 'م', 'گیری', 'تخصص',  
'دارم', '!', 'کلمات', 'عربی', 'مانند', 'اصلاح', 'کاف', 'و',  
'یا\200c\ی', 'برای', 'توکنایزر', 'ما', 'اهمیت', 'دارند', '.',  
'ما', 'می\200c\دانیم', 'که', 'در', 'تاریخ', '2', 'NUM',  
'سپتامبر', '2', 'NUM', '(', '4', 'NUM', 'شهریور', '،', '،', '،', 'NUM'
```

'4' , ' ) ' , ' مقاله\۲00cهای ' , ' » ' , ' ویکی ' , ' پدیا ' , ' « ' , ' ' در ' ,  
, '3' , ' NUM ' , ' زبان ' , ' به ' , ' یک\۲00cمیلیون ' , ' رسید ' , ' . ' , ' ' که ' ,  
, ' این ' , ' مقالات ' , ' شامل ' , ' زمان\۲00cهای ' , ' پیشین ' , ' نیستند ' ,  
, ' . ' , ' در ' , ' ویکی\۲00cپدیا ' , ' فارسی ' , ' ممکن ' , ' است ' , ' ( ' , ' )  
, ' گاهی ' , ' ( ' , ' فاصله ' , ' پراتر\۲00cها ' , ' رعایت ' , ' نشده\_باشد ' ,  
, ' , ' یا ' , ' حتی ' , ' ممکن ' , ' است ' , ' درباره ' ,  
, ' محبوب\۲00cترین\۲00cهای ' , ' فارسی ' , ' صحبت ' , ' شده\_باشد ' ,  
, ' . ' , ' در ' , ' اینجا ' , ' یک ' , ' ایمیل ' , ' آزمایشی ' , ' از ' , ' من ' ,  
, ' , ' ir ' , ' . ' , ' ac ' , ' . ' , ' sh@sbu ' , ' قرار ' , ' دارد ' , ' . ' , ' برای ' ,  
, ' اطلاعات ' , ' بیشتر ' , ' می\۲00ctوانید ' , ' به ' , ' وسایت '  
http' , ' : ' , ' / ' , ' ' آدرس ' , ' به ' , ' ویکی\۲00cpیدا '  
, ' , ' com ' , ' . ' , ' wikipedia ' , ' / ' , ' ' سر ' , ' بزنید ' , ' . ' , ' ) '  
, ' داخل ' , ' پراتر ' , ' بگویم ' , ' , ' , ' این ' , ' یک ' , ' متن ' , ' تستی ' ,  
, ' است ' , ' . ' , ' حداقل ' , ' به ' , ' من ' , ' اینطور ' , ' گفته\۲00cشده '  
, ' است ' , ' . ' , ' ( ' , ' مجله ' , ' تایم ' , ' در ' , ' گزارش ' , ' سال '  
, ' , '4' , ' NUM ' , ' خود ' , ' , ' , ' جیمی ' , ' ویلز ' , ' را ' , ' در ' , ' گروه '  
, ' , '3' , ' NUM ' , ' فرد ' , ' تأثیرگذار ' , ' سال ' , ' اعلام ' , ' کرد ' , ' . ' ,  
, ' همچنین ' , ' در ' , ' همین ' , ' سال ' , ' ویکی ' , ' پدیای ' , ' روسی '  
, ' برنده ' , ' جایزه ' , ' رانت ' , ' ( ' , ' ) ' , ' روسی ' , ' : ' , ' , ' Премия '  
, ' , ' Рунета ' , ' ( ' , ' ) ' , ' در ' , ' بخش ' , ' » ' , ' دانش ' , ' و ' , ' آموزش '  
, ' , ' « ' , ' شد ' , ' . ' , ' این ' , ' جایزه ' , ' از ' , ' طرف ' , ' دولت ' , ' اعطا '  
, ' می\۲00cشود ' , ' . ' , ' همچنین ' , ' ویکی ' , ' پدیا ' , ' جایزه '  
, ' یک\۲00cmیلیون ' , ' دلاری ' , ' مدیریت ' , ' پروژه ' , ' را ' , ' از '  
, ' همایش ' , ' صفاجو ' , ' دریافت ' , ' کرد ' , ' . ' , ' پلنفورم ' , ' اهداف '  
, ' توسعه ' , ' پایدار ' , ' Nations ' , ' United ' , ' : ' , ' چندین ' , ' پروژه '  
, ' متن\_آزاد ' , ' دارد ' , ' که ' , ' وظایف ' , ' غیردانشنامه\۲00cای ' , ' را '  
, ' انجام ' , ' می\۲00cdهند '

## ii. توکن‌بندی با یارسیوار

[ 'عنوان', 'مقاله', ':', 'صفحه', 'اصلی', '>', 'p', '<', 'متن', 'ویکی\200cپدیا', 'فارسی', '/>', 'p', '<', 'من', 'john', 'هستم', 'و', 'در', 'ب', '.', 'م', '.', 'م\200cگیری', 'تخصص', 'دارم', '!', 'کلمات', 'عربی', 'مانند', 'اصلاح', 'کاف', 'و', 'یا', 'ی', 'برای', 'توکنایزر', 'ما', 'اهمیت', 'دارند', '.', 'ما', 'می\200c\200cدانیم', 'که', 'در', 'تاریخ', '20', 'سپتامبر', '2004', '(', '1383', 'y0m6d29', ')', 'مقاله\200cهای', 'ویکی', 'پدیا', 'در', '105', 'زبان', 'به', '1000000', 'رسید', '.', 'که', 'این', 'مقالات', 'شامل', 'زمان\200cهای', 'پیشین', 'نیستند', '.', 'در', 'ویکی\200cپدیا', 'فارسی', 'ممکن', 'است', '(', 'گاهی', ')', 'فاصله', 'پراتنز\200cها', 'رعایت\200cنشده\200cباشد', 'یا', 'حتی', 'ممکن', 'است', 'درباره', 'محبوب', 'ترین\200cهای', 'فارسی', 'صحبت\200cشده\200cباشد', '.', 'در', 'اینجا', '1', 'ایمیل', 'آزمایشی', 'از', 'من', 'sh@sbu', 'ir', 'ac', '.', 'قرار', 'دارد', '.', 'برای', 'اطلاعات', 'بیشتر', 'می\200cتوانید', 'به', 'وسایت', 'ویکی\200cپدیا', 'فارسی', 'به', 'آدرس', 'http://', 'wikipedia', 'com', 'سر', 'بزنید', '.', 'داخل', 'پراتنز', 'بگویم', ']

این '1'، 'متن'، 'تستی'، 'است'، '،'، 'حداقل'، 'به'، 'من'، 'اینطور'، 'گفته' u200c\شده u200c\است'، '،'، 'مجله'، 'تایم'، 'در'، 'گزارش'، 'سال'، 'y0m0d2006'، 'خود'، '،'، 'جیمی'، 'ویلز'، 'را'، 'در'، 'گروه'، '100'، 'فرد'، 'تاثیرگذار'، 'سال'، 'اعلام'، 'کرد'، '،'، 'همچنین'، 'در'، 'همین'، 'سال'، 'ویکی'، 'پدیای'، 'روسی'، 'برنده'، 'جایزه'، 'رانت'، '،'، 'روسی'، '،'، 'Премия Рунета'، '،'، 'در'، 'بخش'، '»'، 'دانش'، 'و'، 'آموزش'، '«'، 'شد'، '،'، 'این'، 'جایزه'، 'از'، 'طرف'، 'دولت'، 'اعطا'، 'می' u200c\شود'، '،'، 'همچنین'، 'ویکی'، 'پدیا'، 'جایزه'، '1000000'، 'دلاری'، 'مدیریت'، 'پروژه'، 'را'، 'از'، 'همایش'، 'صفاجو'، 'دریافت'، 'کرد'، '،'، '،'، 'پلتفورم'، 'اهداف'، 'توسعه'، 'پایدار'، 'United'، 'Nations'، '،'، 'چندین'، 'پروژه'، 'متن'، '–'، 'آزاد'، 'دارد'، 'که'، 'وظایف'، 'غیردانشنامه' u200c\ای'، 'را'، 'انجام'، 'می' u200c\دهند']

## مقایسه

\* امکان استفاده از توکنایزر دادما وجود نداشت.

همانطور که در مرحله نرمال سازی پارسیوار از هضم بهتر عمل کرده است، در این مرحله نیز بهتر عمل می کند.

در هضم به دلیل فاصله نگذاشتن قبل از علائم نگارشی، برخی کلمه‌ها همراه با علائم به عنوان یک توکن در نظر گرفته شده‌اند.

در هضم، در مرحله ی توکن‌بندی تبدیل ایمیل و اعداد به توکن‌های مخصوص را داریم، که با توجه به فاصله‌گذاری‌های نامناسب در مرحله‌ی قبل، قادر به تشخیص ایمیل و url نمی‌باشد.

هر دو توکنایزر کلماتی را که با نیمفاصله به هم متصل شده‌اند را یک توکن در نظر می‌گیرند، که باتوجه به قوی‌تر بودن پارسیوار در این موضوع در مرحله قبل، در این مرحله این پارسیوار متن را به صورت درست‌تری توکنایز کرده‌است.

### c. ریشه‌یابی و بن‌یابی

ریشه‌یابی به معنای حذف بخش‌های افزوده شده به کلمه و تبدیل آن به حالت ساده می‌باشد. دایره لغات مدل حالت منسجم تری داشته باشد. این عمل برای زبان انگلیسی دارای الگوریتمی به نام **porter** می‌باشد، اما در کتابخانه‌های فارسی به صورت کاملاً ابتدایی تعریف شده است. خروجی ریشه‌یابی گاهی کلمات بامعنی نیستند؛ به این دلیل که حذف پسوند تنها با یک شرط انجام می‌شود و در این عملیات کلمه‌ای مثل عالی نیز تحت تاثیر ریشه‌یاب به عال تبدیل می‌شود. (مثال: کتاب‌هایشان -> کتاب)

در فرآیند بن‌یابی بن‌یاب به دنبال یک کلمه با معنی می‌باشد. تفاوت ریشه‌یاب و بن‌یاب در همین مسئله است که کلمات تولید شده بعد از stem کردن ممکن است بامعنی نباشند اما بعد از lemmatize کردن همیشه کلمات با معنی داریم. Lemmatizer های پیشرفته بر اساس morphological parsing پیاده‌سازی شده‌اند، اما کتابخانه های فارسی براساس لیستی از کلمات مرجع به همراه ریشه آن این کار را انجام می‌دهد.

## i.

### ریشه‌یابی با هضم

[ 'عنو' , 'مقاله' , ':' , 'صفحه' , 'اصل' , '<p>متن' ,  
 'ویکی\۲۰۰c\پدیا' , 'فارسی' , '>' , '/' , '<p' , 'من' , 'john' , 'هس' ,  
 'و' , 'در' , 'ب' , '.' , '4' , 'NUM' , 'گیر' , 'تخصص' , 'دار' ,  
 '!' , 'کل' , 'عرب' , 'مانند' , 'اصلاح' , 'کاف' , 'و' , 'یا' , 'برا' ,  
 'توکنایزر' , 'ما' , 'اهم' , 'دارند' , '.' , 'ما' , 'می\۲۰۰c\دان' , 'که' ,  
 'در' , 'تاریخ' , '2' , 'NUM' , 'سپتامبر' , 'NUM' , '4' , 'NUM' ,  
 '2' , 'شهریور' , '4' , 'NUM' , 'مقاله' , '»' , 'ویک' ,  
 'پدیا' , '«' , 'در' , '3' , 'NUM' , 'زب' , 'به' , 'یک\۲۰۰c\میلیون' ,  
 'رسید' , '.' , 'که' , 'این' , 'مقال' , 'شامل' , 'زمان' , 'پیشین' ,  
 'نیستند' , '.' , 'در' , 'ویکی\۲۰۰c\پدیا' , 'فارسی' , 'ممکن' , 'اس' ,  
 ' )' , 'گاه' , ' )' , 'فاصله' , 'پراتنز' , 'رعا' , 'نشده\_باشد' , 'یا' ,  
 'حت' , 'ممکن' , 'اس' , 'درباره' , 'محبوب\۲۰۰c\ترین' , 'فارسی' ,  
 'صحب' , 'شده\_باشد' , '.' , 'در' , 'اینجا' , 'یک' , 'ایمیل' , 'آزمایش' ,  
 'از' , 'من' , 'ir' , '.' , 'ac' , '.' , 'sh@sbu' , 'قرار' , 'دارد' ,  
 ' )' , 'برا' , 'اطلاع' , 'ب' , 'می\۲۰۰c\توانید' , 'به' , 'وبسا' ,  
 'ویکی\۲۰۰c\پیدا' , 'فارسی' , 'به' , 'آدرس' , '/' , ':' , 'http' ,  
 ' )' , 'com' , 'wikipedia' , '/' , 'سر' , 'بزنید' , '.' , ' )' ,  
 'داخل' , 'پراتنز' , 'بگو' , ' )' , 'این' , 'یک' , 'متن' , 'تست' , 'اس' ,  
 ' )' , 'حداقل' , 'به' , 'من' , 'اینطور' , 'گفته\۲۰۰c\شده' , 'اس' ,  
 ' )' , ' )' , 'مجله' , 'تا' , 'در' , 'گزار' , 'سال' , '4' , 'NUM' ,  
 'خود' , ' )' , 'جیم' , 'ویلز' , 'را' , 'در' , 'گروه' , '3' , 'NUM' ,  
 'فرد' , 'تأثیرگذار' , 'سال' , 'اعلا' , 'کرد' , ' )' , 'همچنین' , 'در' ,  
 'همین' , 'سال' , 'ویک' , 'پدیا' , 'روس' , 'برنده' , 'جایزه' , 'ران' ,  
 ' )' , 'روس' , ':' , 'Рунета' , 'Премия' , ' )' , 'در' , 'بخ' ,  
 '«' , 'دان' , 'و' , 'آموز' , '«' , 'شد' , ' )' , 'این' , 'جایزه' , 'از' ,  
 'طرف' , 'دول' , 'اعطا' , 'می\۲۰۰c\شود' , ' )' , 'همچنین' , 'ویک' ,  
 'پدیا' , 'جایزه' , 'یک\۲۰۰c\میلیون' , 'دلار' , 'مدیر' , 'پروژه' , 'را' ,  
 'از' , 'هما' , 'صفاجو' , 'دریاف' , 'کرد' , ' )' , 'پلتفور' , 'اهداف' ,  
 'توسعه' , 'پایدار' , 'United' , 'Nations' , ':' , 'چندین' , 'پروژه' ,  
 'متن-آزاد' , 'دارد' , 'که' , 'وظایف' , 'غیردانشنامه' , 'را' , 'انجا' ,  
 'می\۲۰۰c\دهند' ]

## ii.

### بن‌یابی با هضم

[ 'عنو' , 'مقاله' , ':' , 'صفحه' , 'اصل' , '<p>متن' ,  
 'ویکی\۲۰۰c\پدیا' , 'فارسی' , '>' , '/' , '<p' , 'من' , 'john' , 'هس' ,  
 'و' , 'در' , 'ب' , '.' , '4' , 'NUM' , '#هست' , 'گیر' , 'تخصص' ,  
 'دار' , '!' , 'کل' , 'عرب' , 'مانند' , 'اصلاح' , 'کاف' , 'و' , 'یا' ,  
 'برا' , 'توکنایزر' , 'ما' , 'اهم' , 'داشت#دار' , '.' , 'ما' ,  
 'می\۲۰۰c\دان' , 'که' , 'در' , 'تاریخ' , '2' , 'NUM' , 'سپتامبر' ,  
 '2' , 'NUM' , ' )' , '4' , 'NUM' , 'شهریور' , '4' , 'NUM' ,  
 ' )' , 'مقاله' , '»' , 'ویک' , 'پدیا' , '«' , 'در' , '3' , 'NUM' ,  
 'زب' , 'به' , 'یک\۲۰۰c\میلیون' , 'رسید' , '.' , 'که' , 'این' , 'مقال' ,  
 'شامل' , 'زمان' , 'پیشین' , 'نیستند' , '.' , 'در' , 'ویکی\۲۰۰c\پدیا' ,  
 'فارسی' , 'ممکن' , 'اس' , ' )' , 'گاه' , ' )' , 'فاصله' , 'پراتنز' ,  
 'رعا' , 'شد#شو' , ' )' , 'یا' , 'حت' , 'ممکن' , 'اس' , 'درباره' ,

'محبوب'، 'فارس'، 'صحب'، 'شد#شو'، '،'، 'در'، 'اینجا'، 'یک'،  
 'ایمیل'، 'آزمایش'، 'از'، 'من'، '،'، '،'، 'ac'، '،'، 'sh@sbu'،  
 'ir'، 'قرار'، 'داشت#دار'، '،'، 'برای'، 'اطلاع'، 'ب'،  
 'توانست#توان'، 'به'، 'وبسا'، 'ویکی\200cپیدا'، 'فارس'، 'به'،  
 'آدرس'، 'com'، '،'، 'wikipedia'، '،'، 'http'، '،'، '،'، '،'،  
 'سر'، 'زد#زن'، '،'، '،'، 'داخل'، 'پراتنز'، 'بگو'، '،'، 'این'،  
 'یک'، 'متن'، 'تست'، 'اس'، '،'، 'حداقل'، 'به'، 'من'،  
 'اینطور'، 'گفته\200cشده'، 'اس'، '،'، '،'، 'مجله'، 'تا'، 'در'،  
 'گزار'، 'سال'، '4'، 'NUM'، 'خود'، '،'، 'جیم'، 'ویلز'، 'را'،  
 'در'، 'گروه'، '3'، 'NUM'، 'فرد'، 'تأثیرگذار'، 'سال'، 'اعلا'،  
 'کرد#کن'، '،'، 'همچنین'، 'در'، 'همین'، 'سال'، 'ویک'، 'پدیا'،  
 'روس'، 'برنده'، 'جایزه'، 'ران'، '،'، 'روس'، '،'، '،'، 'Премия'،  
 'Рунета'، '،'، 'در'، 'بخ'، '»'، 'دان'، 'و'، 'آموز'، '«'،  
 'شد#شو'، '،'، 'این'، 'جایزه'، 'از'، 'طرف'، 'دول'، 'اعطا'،  
 'شد#شو'، '،'، 'همچنین'، 'ویک'، 'پدیا'، 'جایزه'،  
 'یک\200cمیلیون'، 'دلار'، 'مدیر'، 'پروژه'، 'را'، 'از'، 'هما'،  
 'صفاجو'، 'دریاف'، 'کرد#کن'، '،'، 'پلتفور'، 'اهداف'، 'توسعه'،  
 'پایدار'، 'United'، 'Nations'، '،'، 'چندین'، 'پروژه'،  
 'متن-آزاد'، 'داشت#دار'، 'که'، 'وظایف'، 'غیردانشنامه'، 'را'، 'انجا'،  
 'داد#ده' ]

### iii.

### ریشه و بنیابی با پارسیوار

[ 'عنوان'، 'مقاله'، '،'، '،'، 'اصلی'، '،'، '>'، 'p'، '،'، '<'، 'متن'،  
 'ویکی\200cپدیا'، 'فارسی'، '،'، '<'، 'p'، '،'، '>'، 'john'،  
 'هست'، 'و'، 'در'، 'ب'، '،'، 'م'، '،'، 'م\200cگیری'،  
 'تخصص'، 'داشت&دار'، '!'، 'کلمات'، 'عربی'، 'مانند'، 'اصلاح'،  
 'کاف'، 'و'، 'یا'، 'ی'، 'برای'، 'توکنایزر'، 'ما'، 'اهمیت'،  
 'داشت&دار'، '،'، 'ما'، 'دانست&دان'، 'که'، 'در'، 'تاریخ'، '20'،  
 'سپتامبر'، '2004'، '،'، '1383'، '،'، 'y0m6d29'، '،'، 'مقاله'،  
 '،'، 'ویکی'، 'پدیا'، '،'، 'در'، '105'، 'زبان'، 'به'،  
 '1000000'، 'رسید'، '،'، 'که'، 'این'، 'مقالات'، 'شامل'، 'زمان'،  
 'پیشین'، 'نیستند'، '،'، 'در'، 'ویکی\200cپدیا'، 'فارسی'، 'ممکن'،  
 'اس'، '،'، 'گاهی'، '،'، 'فاصله'، 'پراتنز'،  
 'رعایت\200cنشده\200cباشد'، '،'، 'یا'، 'حتی'، 'ممکن'، 'اس'،  
 'درباره'، 'محبوب'، 'ترین\200cهای'، 'فارسی'،  
 'صحت\200cشده\200cباشد'، '،'، 'در'، 'اینجا'، '1'، 'ایمیل'،  
 'sh@sbu'، '،'، 'ac'، '،'، 'ir'، 'من'،  
 'قرار'، 'داشت&دارد'، '،'، 'برای'، 'اطلاعات'، 'بیشتر'،  
 'توانست&توان'، 'به'، 'وبسایت'، 'ویکی\200cپیدا'، 'فارسی'، 'به'،  
 'آدرس'، 'com'، '،'، 'wikipedia'، '،'، 'http'، '،'، 'سر'،  
 'زد&زن'، '،'، '،'، 'داخل'، 'پراتنز'، 'گفت&گو'، '،'، 'این'، '1'،  
 'متن'، 'تست'، 'اس'، '،'، 'حداقل'، 'به'، 'من'، 'اینطور'،  
 'گفته\200cشده\200cاست'، '،'، 'مجله'، 'تایم'، 'در'، 'گزارش'،  
 'سال'، 'y0m0d2006'، 'خود'، '،'، 'جیمی'، 'ویلز'، 'را'، 'در'،  
 'گروه'، '100'، 'فرد'، 'تأثیرگذار'، 'سال'، 'اعلام'، 'کرد'، '،'،  
 'همچنین'، 'در'، 'همین'، 'سال'، 'ویکی'، 'پدیای'، 'روسی'،  
 'برنده'، 'جایزه'، 'ران'، '،'، 'روسی'، '،'، '،'، 'Премия'،  
 'Рунета'، '،'، 'در'، 'بخش'، '»'، 'دانش'، 'و'، 'آموزش'،



'«', 'شد', '!', 'این', 'جایزه', 'از', 'طرف', 'دولت', 'اعطا', 'شد', 'شو', '!', 'همچنین', 'وبی', 'پدیا', 'جایزه', '1000000', 'دلاری', 'مدیریت', 'پروژه', 'را', 'از', 'همایش', 'صفاجو', 'دریافت', 'کرد', '!', 'پلتفرم', 'اهداف', 'توسعه', 'پایدار', 'United', 'Nations', ':', 'چندین', 'پروژه', 'متن', '-', 'آزاد', 'داشت&دارد', 'که', 'وظایف', 'غیردانشنامه\200cای', 'را', 'انجام', 'داد&ده']

#### iv. مقایسه

\* پارسیوار دارای بن‌یاب و ریشه‌یاب جداگانه نیست و این کار را به صورت ترکیبی در FindStems انجام می‌دهد. در هضم stemmer اگر کلمات با ["ات", "ان", "ترین", "تر", "م", "ت", "ش", "یی", "ی", "ها", "ا", ""] تمام شده باشند، آن بخش‌ها را حذف می‌کند و ه را به ه تبدیل می‌کند. به همین دلیل دارای خطای بسیار است. به طور مثال کلمه ی اعلام را به اعلا تبدیل می‌کند. Lemmatizer هضم ابتدا کلمه را با لیستی از کلمات ذخیره شده مقایسه می‌کند، در صورت مچ شدن همین کلمه را برمی‌گرداند. در غیر این صورت، در فایلی از افعال به دنبال کلمه می‌گردد تا ریشه را بازگرداند. سپس یک بار دیگر کلمه را stem می‌کند و با لغات ذخیره شده چک می‌کند. در نهایت یا stem شده را برمی‌گرداند یا خود کلمه را. Stemmizer پارسیوار خیلی ضعیف عمل می‌کند. به خصوص برای کلماتی که دارای نیم‌فاصله می‌باشند. در مراحل مختلف پسوند و پیشوندهای افعال و دیگر کلمات را حذف می‌کند. تمرکز آن بیشتر روی افعال می‌باشد و توانایی ریشه‌یابی کلمات عادی آن بسیار پایین است؛ به طور مثال نمی‌تواند غیردانشنامه‌ای را ساده سازی کند. همچنین در بن‌یابی فعل ساده‌ای مثل هستند شکست می‌خورد.

### 3. نتیجه‌گیری

پکیج دادما دارای باگ‌های بسیاری بود و کارکردن با آن بسیار سخت! پارسیوار و هضم دارای تفاوت‌هایی در مراحل مختلف بودند که باید در هنگام انتخاب به آنها توجه کرد و با توجه به نوع متن و همچنین تسک مورد نظر یک کدام را برگزید. به طور کلی پارسیوار در مرحله توکنایز کردن بهتر عمل کرد اما هضم برای ریشه‌یابی و بن‌یابی بهتر عمل می‌کند.

### 4. منابع

<https://github.com/roshan-research/hazm/tree/master>

<https://github.com/Dadmata/DadmaTools/tree/main>

<https://github.com/ICTRC/Parsivar/tree/master>  
<https://web.stanford.edu/~jurafsky/slp3/ed3book.pdf>