

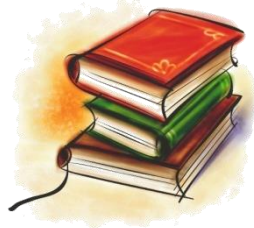
آشنایی با زبان‌شناسی رایانشی

واحدسازی و نرمال‌سازی

هادی ویسی

h.veisi@ut.ac.ir

دانشگاه تهران - دانشکده علوم و فنون نوین



فهرست

○ تعریف واحدسازی و نرمال‌سازی

○ واحدسازی

- مشکلات
- راه‌حل‌ها
- مشکلات زبان فارسی و راه‌حل‌ها

○ نرمال‌سازی

- مشکلات زبان فارسی
- راه‌حل‌ها
- موارد تکمیلی



واحدسازی و نرمال‌سازی

- برای پردازش متن، اولین کار واحدسازی و نرمال‌سازی متن ورودی است
 - کاربردهایی مانند ترجمه ماشینی، بازیابی اطلاعات، تبدیل متن به گفتار و ...

- واحدسازی (Tokenization)
 - تقطیع دنباله کاراکترهای تشکیل دهنده متن به دنباله ای از کلمات

- نرمال‌سازی (Normalization)
 - یکدست‌سازی واحدهای متنی به طوری که قابل پردازش توسط ماشین باشد.



واحدسازی

○ ساده‌ترین راه برای واحدسازی متن و تقطیع آن به دنباله کلمات

- واحدسازی بر اساس کاراکتر فاصله (Space)
- Token: واحدهای متنی که با استفاده از فاصله از هم جدا می‌گردند.

○ مشکل: همیشه کاراکتر فاصله مرز کلمات را مشخص نمی‌کند

- کلماتی که از چند token تشکیل شده اند (Multi Token Units)

○ می‌روند - کتاب‌ها - در حالی که - زبان‌شناسی

○ air force - St. Petersburg - go back - take off

- کلماتی که به هم چسبیده و یک Token را به وجود آورده‌اند (Multi Unit Tokens)

○ ناشی از حروفی که شکل متصل و غیرمتصل آنها با هم یکسان است: ر، ژ، و، ذ، د، آ، ا

○ در کتاب ⇐ در کتاب

○ مرا/زوی/کزین ⇐ من/را/از وی/که از این (برای POS tagging)

○ و با ⇐ و با («وبا»؟)



واحدسازی: مشکلات

○ علائم نشانه‌گذاری (Punctuations) معمولاً به کلمه قبل از خود می‌چسبند.

- علائم نقطه گذاری را هم مانند کاراکتر فاصله به عنوان مرز کلمات به حساب آوریم.
- استثنا: اختصاراتی مانند B.B.C و اعداد ممیزدار مانند 2.3

○ پسوندها و پیشوندها

- در فارسی پسوندها و پیشوندها در بسیاری مواقع با فاصله از کلمه اصلی نوشته می‌شوند
- مثال: رفته ام – بر می‌گردد – دانش آموز

○ کلمات مرکب

- بین اجزای کلمات مرکب معمولاً فاصله گذاشته می‌شود
- دوچرخه سوار – برون مرزی



واحدسازی: راه حل‌ها

○ استفاده از علائم نشانه‌گذاری

- تعیین قطعی مرز کلمات با علائم «،»، «:»، «!»، «؟»
- نقطه «.»: علاوه بر مرز در سرنام‌ها و علائم اختصار هم می‌آید

○ استفاده از قواعد املائی برای پسوندها و پیشوندها

- مثال: چسباندن تمام پیشوندهای "می" به اول کلمات و پسوندهای "ها" به آخر کلمات
- رفع موارد استثنا: مشخص کردن تمام کلمات مجازی که وندها به آنها می‌چسبند (با استفاده از یک پیکره متنی بزرگ)

○ استفاده از یک پیکره متنی بزرگ برای مشخص نمودن MUT‌های پرکاربرد و تشکیل یک جدول از آنها برای تقطیع

- استفاده از یک پیکره متنی تقطیع شده برای مشخص کردن کلمات چندقسمتی و تشکیل قواعد املائی برای آنها
- استفاده از جدول حاوی لیست کلمات مرکب



نرمال‌سازی: مشکلات زبان فارسی ...

○ وجود encoding‌های مختلف برای بعضی از کاراکترها
• مانند “ی” و “ک”

○ روش‌های مختلف چسبیدن وندها به کلمات اصلی
• می روند - می‌روند - میروند
• زبان شناسی - زبان‌شناسی - زبانشناسی

○ روش‌های مختلف اتصال اجزای کلمات مرکب
• برون مرزی - برون‌مرزی - برونمرزی

○ کلمات چنداملایی
• مسئولیت - مسوولیت - مسؤولیت
• گفتگو - گفت و گو - گفت وگو - گفت‌وگو
• پاییز - پائیز



نرمال‌سازی: مشکلات زبان فارسی ...

○ فرم‌های نوشتاری مختلف کلمه «درحالی‌که» در پیکره متنی فارسی

تعداد	صورت نوشتاری
22,845	در_حالی_که
2,210	درحالی^که
2,057	درحالی_که
1,818	در_حالیکه
1,220	در_حالی^که
359	درحالیکه
24	در^حالی^که
4	در^حالی_که
2	در^حالیکه

○ ابهام در تبدیل فاصله به نیم‌فاصله

- اسب سواری یکی از تفریحات انسان‌هاست
- او اسب سواری را با تازیانه زد
- من با اسب سواری کردم



نرمال‌سازی: راه حل‌ها ...

○ درست کردن فاصله بین متن و علائم نشانه‌گذاری

- علائم «.»، «،»، «؛»، «:»، «!» و «؟» به کلمات قبلی می‌چسبند و با کلمه بعدی به اندازه یک واحد فاصله (Space)، فاصله دارند

○ مثال: « من، تو ، او ضمائر شخصی فارسی هستند . » ⇨ « من، تو، او ضمائر شخصی فارسی هستند. »

○ <عبارت+ علامت+ فاصله+ عبارت>

○ استثنا: علامت «.» در سرنام‌ها و اختصارها

○ مثال: « در سال پنجم ه. ق » ⇨ « در سال پنجم ه.ق »

- علائم { } ، [] ، () و « » به عبارت داخل خود می‌چسبند

○ مثال: « دانشجویان اینجا (دانشگاه تهران) در ... » ⇨ « دانشجویان اینجا (دانشگاه تهران) در ... »

○ <عبارت+ فاصله+ علامت ابتدا+ عبارت داخل + علامت انتهایی+ فاصله+ عبارت>



نرمال‌سازی: راه حل‌ها ...

○ درست کردن فاصله بین متن و اعداد

- اعداد باید با کلمات قبل و بعد از خود یک واحد فاصله داشته باشند

○ مثال: «دیروز ۶مرد و ۴زن مدل گرفتند.» ⇨ «دیروز ۶مرد و ۴زن مدل گرفتند»

- در اعداد اعشاری و تاریخ علامت «ممیز» باید یک‌دست و اصلاح شود

○ مثال: «معدلش در ۱۳۹۳/۱۱/۱۲ برابر با ۱۷.۵ است.» ⇨ «معدلش در ۱۳۹۳/۱۱/۱۲ برابر با ۱۷/۵ است.»



نرمال‌سازی: راه حل‌ها ...

○ درست کردن فاصله بین کلمات و پسوندها

- پسوندها: اشتقاقی یا صرفی
- < کلمه + پسوند اشتقاقی + پسوند صفت + نشانه جمع + پسوند ملکی (یا ی نکره) >
 - پسوند صفت: مانند «تر» و «ترین»
 - نشانه جمع: مانند «ها»
 - پسوند ملکی: مانند «م، ت، ش، مان، تان، شان»
 - یک کلمه نمی‌تواند هم شامل «ی نکره» و هم پسوند ملکی باشد
- کلماتی که به «ا» یا «و» ختم می‌شوند، قبل از پسوند ملکی یک «ی» میانجی قرار می‌گیرد
 - مثال: عصایم، سبویم
- برخی پسوندهای اشتقاقی می‌توانند نقش اسم هم داشته باشند
 - مثال: «سیر» در «گرم سیر» یا «سیر ≠ گرسنه»
- کلماتی که چند پسوند دارند، همه آنها با نیم‌فاصله از کلمه (یا سایر پسوندها) جدا می‌شوند
 - مثال: «جزوه‌ها یمان» ⇐ «جزوه‌هایمان»
 - مثال: «حیله‌گرها» ⇐ «حیله‌گرها»
 - مثال: «زشت‌ترین‌ها» ⇐ «زشت‌ترین‌ها»



نرمال‌سازی: راه حل‌ها ...

○ درست کردن فاصله بین کلمات و پیشوندهای اشتقاقی

• <پیشوند اشتقاقی + کلمه>: با نیم‌فاصله از کلمه بعد قرار می‌گیرند

○ مثال: « بی ادب » ⇐ « بی‌ادب »

○ مثال: « نا مرد » ⇐ « نامرد »



نرمال‌سازی: راه حل‌ها ...

○ درست کردن فاصله در افعال

- پسوندهای فعل ماضی نقلی: با نیم‌فاصله از کلمه قبل (فعل) قرار می‌گیرند
 - همه پسوندهای «ام»، «ای»، «است»، «ایم»، «اید» و «اند»
 - گاهی «است» را استثنا می‌کنند و آن را با فاصله جدا می‌کنند
 - مثال: « گرفته ام » ⇐ « گرفته‌ام »
- پسوند «می» در ماضی استمراری و مضارع اخباری: با نیم‌فاصله از جز فعلی قرار می‌گیرد
 - استثنا: «می به معنی شراب» ⇐ اگر کلمه بعد از «می» اسم بود، فاصله را حذف نکنید
 - قطعی نیست: «می» در «می خورد» می‌تواند هم به معنای شراب باشد و هم جز فعلی
- اصلاح پسوندهایی مانند «بر»
 - مثال: « بر خواهد گشت » ⇐ « برخواهد گشت »



نرمال‌سازی: راه حل‌ها ...

○ درست کردن نشانه جمع «ها»

- هم به صورت متصل و هم به صورت مجزا ظاهر می‌شود
- مواردی که «ها» باید جدا نوشته شود
 - قرار گرفتن بعد از کلمات بیگانه نامانوس مانند «فرمالیست‌ها»
 - مشخص کردن کلمه به منظور برجسته کردن آن و یا آموزش: « دبستانی‌ها»
- در نرمال‌سازی بهتر است همه جا «ها» را با نیم‌فاصله جدا کنید



نرمال‌سازی: راه‌حل‌ها ...

○ یک‌دست‌سازی املاي کلمات چنداملايي با استفاده از جدول حاوي کلمات (Lookup Table) بر اساس شیوه نگارش استاندارد

- مسوولیت ← مسئولیت

○ حذف علائم نشانه‌گذاری (diacritics)

- فناوری ← فناوری
- کاملاً ← کاملاً

○ درست کردن کاراکترهایی که چند کد دارند

- ی ← ی
- ک ← ک
- و ← و
- ی ← ی
- ا ← ا
- ه ← ه



نرمال‌سازی: راه‌حل‌ها

○ سایر

- حذف «ی» از آخر کلماتی مانند «کلمه‌ی»
- حذف «ء» از کلماتی مانند «املاء»
- چسپاندن پیشوند «هم» در کلماتی مانند «هم‌چنین»
- چسپاندن پیشوند «این» و «آن» در کلماتی مانند «این‌جا» و «آن‌جا»

○ نیم فاصله

- کاراکتر zero-width non-joiner (ZWNJ)
- یونی‌کد: U+200C
- در HTML: `‌` یا `&zwnj`



نرمال‌سازی: موارد تکمیلی

○ در بعضی کاربردها مانند تبدیل متن به گفتار، موارد دیگری نیز باید نرمال‌سازی شوند

• بعضی از اختصارات

○ (ع) ← علیه السلام

• اعداد

○ ۲۳ ← بیست و سه

○ ۲.۳ ← دو ممیز سه دهم

○ ۱/۱۰۰۰ ← یک هزارم

○ ۶۶۶۱۱۵۲۵

○ اگر شماره تلفن باشد: شصت و شش شصت و یک پانزده بیست و پنج

○ اگر مبلغ باشد: شصت و شش میلیون و شش صد و یازده هزار و پانصد و بیست و پنج

• تاریخ

○ ۸۹/۸/۱۵ ← پانزده هشت هشتاد و نه