

توضیحات پیش پردازش داده‌های متنی با استفاده از Regular Expression

عرفان محمدزاده

er.mohammadzadeh@gmail.com

 [e-mohammadzadeh](#)

۱۴۰۴ آذر ۵

فهرست مطالب

سه	فهرست جداول
چهار	فهرست تصاویر
۱	۱ مقدمه
۲	۲ خواندن مجموعه داده
۳	۳ حذف
۲	Remove Null Values ۱.۳
۲	Remove Stopwords ۲.۳
۳	Remove Date ۳.۳
۷	Remove Tag ۴.۳
۷	Remove Cashtag ۵.۳
۸	Remove Number ۶.۳
۹	Remove Title ۷.۳
۱۰	Remove Punctuation ۸.۳
۱۰	۴ جاگجایی
۱۰	Slang Words ۱.۴
۱۱	Lowercase ۲.۴
۱۱	Contraction ۳.۴
۱۲	URL ۴.۴
۱۳	Email Address ۵.۴
۱۳	Currency ۶.۴
۱۵	Time ۷.۴
۱۶	Repeated Letter ۸.۴

یک

۱۶	۵ بسط دادن
۱۶	Mention ۱.۵
۱۷	Hashtag ۲.۵
۱۷	۶ ذخیره مجموعه داده
۱۸	۷ مراجع

فهرست جداول

۱	ترتیب مراحل پیش‌پردازش	۱
۴	انواع فرمتهای کوتاه برای زمان	۲
۵	انواع فرمتهای کامل برای زمان	۳
۶	انواع فرمتهای برای اسمی ماهها	۴
۱۴	علائم اختصاری واحد پول کشورها	۵
۱۵	سیستم شمارش پول	۶
۱۵	سیستم تشخیص ساعت	۷

فهرست تصاویر

لیست ایستوازه های غیرمنفی در زبان انگلیسی	۱
لیست عناوین زبان انگلیسی	۲
بخشی از کلمات عامیانه زبان انگلیسی به همراه شکل کامل هر کدام ..	۳
بخشی از کلمات اختصار زبان انگلیسی به همراه شکل کامل هر کدام ..	۴

۱ مقدمه

برای اینکه بتوانیم مجموعه‌داده‌ها را با مدل تجزیه و تحلیل کنیم، قبل از آن نیاز داریم که این مجموعه‌داده‌ها را آماده‌سازی کنیم. آماده‌سازی یا پیش‌پردازش مجموعه‌داده به پاک‌سازی داده از موارد ناخواسته گفته می‌شود، مثل: پاک کردن کلمات یا علائم اضافی یا فضای خالی از مجموعه‌داده. در ابتدا نحوه خواندن مجموعه‌داده موردنظر را توضیح می‌دهیم و سپس مواردی که برای پیش‌پردازش مجموعه‌داده‌هایمان انجام شده را با جزئیات شرح می‌دهیم. قسمت پیش‌پردازش به سه بخش اصلی تقسیم می‌شود: جابجایی^۱، حذف^۲ و بسط دادن^۳. ترتیب مراحل پیش‌پردازش انجام شده، برابر با ترتیبی است که در جدول ۱ نوشته شده است.

#	Preprocessing Steps
1	Remove Null Values
2	Replace Slang words
3	Lowercase
4	Expand Contractions
5	Remove Stopwords
6	Replace URLs
7	Replace Email Addresses
8	Expand Mentions
9	Expand Hashtags
#	Preprocessing Steps
10	Remove Date
11	Replace Currencies
12	Replace Time
13	Remove Tags
14	Remove Cashtags
15	Remove Numbers
16	Remove Titles
17	Remove Punctuations
18	Replace Repeated Letters

جدول ۱: ترتیب مراحل پیش‌پردازش

¹Replace

²Remove

³Expand

۲ خواندن مجموعه‌داده

در ابتدا آدرس^۱ و نام فایل‌های موردنظر به برنامه داده می‌شود، سپس فایل‌ها به ترتیب خوانده شده و در یک دیکشنری^۲ ذخیره می‌شوند. قسمت کلید^۳ در دیکشنری یکی از موارد train، dev و test می‌تواند باشد. قسمت مقدار^۴ در دیکشنری، فایل خوانده شده متناظر با کلید است. سپس مراحل پیش‌پردازش بر روی هر کدام از مقادیر این دیکشنری، اعمال می‌شود و حاصل یک داده تیزی است که در یک دیکشنری دیگر ذخیره می‌گردد. در ادامه هر کدام از مراحل پیش‌پردازش با جزئیات شرح داده شدند.

در ابتدا، ستون‌های مجموعه‌داده را با توجه به محتوای آن‌ها دو بخش می‌کنیم. یک قسمت ستون‌هایی که داده‌های عددی دارند و یک قسمت ستون‌هایی که داده‌های متنی دارند. سپس مراحل پیش‌پردازش بر روی ستون‌های حاوی داده‌های متنی اعمال می‌شوند.

۳ حذف

Remove Null Values ۱۰۳

بخش اول، مربوط به حذف برخی کاراکترها، کلمات و یا علائم خاص است. اگر مجموعه‌داده‌ای که در این بخش بررسی می‌شود، حاوی سلولی^۵ با مقدار تهی^۶ باشد، آن سطر کلاً پاک می‌شود هر چند که سایر سلول‌های آن سطر حاوی مقادیری باشند. همچنین برای اطلاع کاربر، شماره سطرهایی که سلول تهی دارند و قرار هست پاک شوند، در خروجی چاپ می‌شود.

Remove Stopwords ۲۰۳

یک فایل متنی که توسط نویسنده ساخته شده، در کنار کد قرار دارد که دارای فرمت txt است. این فایل محتوی اکثر ایست‌واژه^۷‌های زبان انگلیسی است با این تفاوت که کلمات منفی را به دلیل اینکه نمی‌خواهیم

¹Path

²Dictionary

³Key

⁴Value

⁵Cell

⁶Not a Number (NaN)

⁷Stopword

English Stopwords.txt

1: until, their, further, can, each, yourself, it, myself, out, were,
2: will, but, where, ve, should've, above, your, again, up, me, those,
3: an, very, these, needn, having, he, under, how, m, between, its,
4: about, had, this, that'll, it's, they, hers, when, any, she, have,
5: of, for, during, we, while, below, the, she's, through, herself,
6: before, if, you've, other, now, that, own, off, ourselves, you're,
7: with, whom, and, has, into, in, on, so, d, most, them, itself, same,
8: down, you'll, is, should, because, from, yours, then, themselves,
9: such, i, over, there, being, or, at, been, her, ours, did, here,
10: a, his, are, you'd, y, just, why, than, yourselves, our, be, which,
11: am, theirs, doing, was, s, ll, after, more, what, re, my, both, do,
12: does, all, o, to, himself, as, you, who, only, by, too, t, once,
13: against, few, ma, him, some

شکل ۱: لیست ایستوازه های غیرمنفی در زبان انگلیسی

کلمات منفی از مجموعه داده حذف شوند در نظر نگرفتیم. این فایل متنی شامل ۱۴۰ ایستوازه انگلیسی است. شکل ۱ لیست ایستوازه های انگلیسی را نشان می دهد با این تفاوت که کلمات منفی از این لیست حذف شدند.

Remove Date ۳۰۳

هدف از این قسمت، استخراج و حذف تاریخ است. نحوه بیان تاریخ می تواند متفاوت باشد، سعی شده است که اکثر فرمتهای مهم برای بیان تاریخ را استخراج کنیم. جدول ۲ فرمتهای کوتاه و جدول ۳ فرمتهای کامل را که برای زمان می توانیم از متن استخراج بکنیم نشان می دهد. انواع حالت هایی که برای اسامی ماه های سال میلادی می توان در نظر گرفت نیز در شکل ۴ ذکر شده است. برای استخراج تاریخ از متن هم می توان به تنها یی از فرمت کوتاه و کامل استفاده کرد و هم می توان این دو فرمت را با هم ترکیب کرد. پس در نتیجه تعداد حالت هایی که می توان استخراج کرد زیاد هست و این باعث دقت بالای برنامه در هنگام پیش پردازش می شود.

جدول ۲ : انواع فرمتهای کوتاه برای زمان

#	Short Date		
1	01/01/2024	01-01-2024	01.01.2024
2	27/5/24	27-5-24	27.5.24
3	01/2024	01-2024	01.2024

جدول ۳: انواع فرمتهای کامل برای زمان

Full Date			
#	First Part	Second Part	Third Part
1	Month-name*	of	Any number {4 digit}
2	Month-name	of	Any number {4 digit},
3	Any number {1-4 digit}	Month-name	
4	Any number {1-4 digit},	Month-name,	
5	Month-name	Any number {2 digit} {st, nd, rd, n-th}	Any number {1-4 digit}
6	Month-name,	Any number {2 digit} {st, nd, rd, n-th},	Any number {1-4 digit}
7	Month-name	Any number {1-4 digit}	Any number {2 digit} {st, nd, rd, n-th}
8	Month-name,	Any number {1-4 digit}	Any number {2 digit} {st, nd, rd, n-th},
9	Any number {1-4 digit}	Month-name	Any number {2 digit} {st, nd, rd, n-th}
10	Any number {1-4 digit},	Month-name,	Any number {2 digit} {st, nd, rd, n-th}

* Each month name (table 4) can be substituted

جدول ۴: انواع فرمتهای برای اسامی ماهها

Month Names			
Jan	January	jan	january
Feb	February	feb	february
Mar	March	mar	march
Apr	April	apr	april
-	May	-	may
Jun	June	jun	june
Jul	July	jul	july
Aug	August	aug	august
Sept	September	sept	september
Oct	October	oct	october
Nov	November	nov	november
Dec	December	dec	december

Remove Tag ۴.۲

در این بخش، تگ^۱ های شروع و پایان HTML و سایر تگ‌ها را شناسایی و حذف می‌کنیم. به عنوان مثال، تگ‌های رو برو در این بخش حذف می‌شوند:

```
1: <head>
2: <url>
3: </body>
4: </div>
5: 
```

Remove Cashtag ۵.۲

طبق تعریف شرکت ایکس^۲ (توبیت سابق)، «برچسب مالی یک نماد منحصر به فرد و مخفف برای یک رمزارز^۳ یا یک شرکتی است که قبل از آن علامت دلار آمریکا قرار دارد» [۱]. از برچسب‌های مالی در برنامه ایکس استفاده می‌شود. حال در این قسمت، قصد داریم که این برچسب‌های مالی را در صورت وجود از متن استخراج و حذف کنیم. به عنوان مثال، برچسب‌های مالی رو برو و مشابه آن‌ها در این بخش حذف می‌شوند:

- 1 \$GOOG → برچسب مالی شرکت گوگل
- 2 \$TSLA → برچسب مالی شرکت تسلا
- 3 \$PEP → برچسب مالی شرکت پپسی
- 4 \$BTC → برچسب مالی رمز ارز بیت‌کوین
- 5 \$ETH → برچسب مالی رمز ارز اتریوم

¹Tag

²X

³Cryptocurrency

Remove Number ۶۰۳

پاک کردن اعداد موجود در متن، یکی دیگه از مراحل پیش‌پردازش است. اعدادی که حذف می‌شوند، چنین فرمتی می‌توانند داشته باشند:

- عدد مثبت یا منفی
- عدد کسری
- عدد اعشاری
- عددی که سه رقم سه رقم با کاراکتر \pm یا ، جدا شده است
- عدد علمی به فرمت $Ae \pm B$ یا $AE \pm B$ هر دو عدد هستند).

هنگام پاک‌سازی متن از اعداد با دو چالش مواجه شدیم که در ادامه به صورت مختصر به هر کدام اشاره‌ای شده است:

- پاک کردن همه اعداد موجود در داخل متن به طوریکه اعداد داخل ساختار

$\triangleright \{L/M/T/E\}_{counter} \triangleleft$

دست نخورده باقی بمانند. برای رفع این مشکل، ابتدا این ساختار را شناسایی کرده و سپس نوشه‌های داخل دو کاراکتر \triangleright و \triangleleft را برش داده و بعد از آن عملیات پاک کردن اعداد را انجام می‌دهیم. و در نهایت نوشه‌های برش داده شده در جای خود (بین دو کاراکتر) الصاق می‌شوند.

- دومین چالش، پاک کردن اعداد به طوریکه اعداد داخل کاراکترهای یونیکد^۱ حذف نشوند. وقتی در مجموعه‌داده جمله یا کلمه‌ای غیر انگلیسی وجود داشته باشد، احتمال وجود کاراکتر یونیکد در آن کلمه یا جمله وجود دارد. کاراکترهای یونیکد معمولاً دارای چنین ساختاری هستند:

\u0000

\u10FFFF

^۱Unicode Characters

بعد از حرف «u» شش کاراکتر می‌توان نوشت که بازه هر کاراکتر «اعداد بین صفر تا نه - حروف کوچک انگلیسی بین a تا f - و همچنین حروف بزرگ انگلیسی بین A تا F» است. در ادامه چندین مثال از کاراکترهای یونیکد نوشته شده است:

K\u00f6lsch → Kölsch

pur\u00e9ed → puréed

jalape\u00f1o → jalapeño

Gar\u00e7on → Garçon

با استفاده از دستوراتی ابتدا کاراکترهای یونیکد را شناسایی کرده و سپس شش کاراکتر را نادیده می‌گیریم و باقی اعداد را حذف می‌کنیم و بدین صورت دومین چالش هنگام حذف اعداد نیز رفع می‌شود.

Remove Title ۷۰۲

حذف عناوین استفاده شده در متن نیز می‌تواند جزو یکی از کارهای پیش‌پردازش باشد، عناوینی مثل

Mr. Mrs. Miss. Dr. Prof. Pres. (president) ...

که به وفور در متن‌های انگلیسی تکرار می‌شوند. هفده کلمه انگلیسی برای این بخش در نظر گرفتیم که لیست کامل کلمات در عکس ۲ قابل مشاهده هستند. برای پاک کردن عناوین این نکات در نظر گرفته شدند: اولاً عدم حساسیت به حروف بزرگ و کوچک؛ ثانياً وجود یا عدم وجود نقطه بعد از عنوان.

English Titles

Mr, Ms, Mrs, Miss, Dr, Prof, Sir, Ma'am, Madam, Madame,
Rev, # Reverend, used for Christian clergy
Fr, # Father, used for Catholic priests
Sr, # Sister, used for Catholic nuns
Capt, # Captain, used in the military and for pilots
Gen, # General, used in the military
Hon, # Honorable, used for judges and certain politicians
Pres, # President, used for presidents of companies or organizations

شكل ۲: لیست عناوین زبان انگلیسی

Remove Punctuation ۸.۳

حذف علائم نگارشی و بصری اضافی از متن نیز از مراحل پیش‌پردازش به حساب می‌آید. در این قسمت، کاراکترهایی از قبیل

‘ - = [] \ ; ’ , . / ! @ # \$ % ^ & * () _ + | : ” < > ?

را از مجموعه‌داده پیدا و حذف می‌کنیم.

در این قسمت نیز چالش مربوط به کاراکترهای یونیکد وجود داشت؛ کاراکتر یونیکد بدون علامت «» معنی ندارد و نمی‌تواند به شکل موردنظر تبدیل شود. بنابراین هنگام حذف علائم نگارشی باید دقت کنیم که علامت «» را در کاراکترهای یونیکد حذف نکنیم.

۴ جابجایی

بخش دوم، مربوط به جابجایی برخی کاراکترها، کلمات یا علائم خاص با کلمات دیگر یا کلیدهای منحصر به فردی است.

Slang Words ۱۰۴

اولین کار در قسمت جابجایی، مربوط به جایگزینی کلمات عامیانه زبان انگلیسی با فرم کامل و رسمی است. در نوشتن متن به زبان انگلیسی، برای اینکه بتوان در سریع‌ترین زمان ممکن عملیات تایپ و ارسال پیام را انجام داد، معمولاً از کلمات مخفف و عامیانه زیاد استفاده می‌شود. برای اینکه بتوان دقت مدل را افزایش داد، ما سعی کردیم تا این کلمات مخفف و عامیانه را تا حدودی شناسایی و با کلمات کامل هر کدام جایگزین کنیم.

برای این منظور یک فایل [۲] که قبل ساخته شده بود را تغییراتی دادیم؛ در نهایت فایل حاوی یک دیکشنری است که key همان کلمه مخفف و value همان کلمه به صورت کامل است. شکل ۳ نمونه‌ای از این فایل را نشان می‌دهد. برای دسترسی به فایل کامل، به گیت‌هاب^۱ نویسنده مراجعه شود.

¹GitHub

English Slang Words.json

```
1: {"07734": "hello",
2: "2day": "today",
3: "2ge4": "Together",
4: "2morrow": "tomorrow",
5: "4ever": "forever",
6: "0noe": "Oh No",
7: "0vr": "over",
8: "10q": "thank you",
9: "5n": "fine",
10: "absnt": "absent",
11: "bc": "because",
12: "c@": "cat",
13: "dw": "don't worry",
14: ...}
```

شکل ۳: بخشی از کلمات عامیانه زبان انگلیسی به همراه شکل کامل هر کدام

Lowercase ۲.۴

یکی از کارهایی که بر روی مجموعه داده‌ها انجام می‌گیرد، تبدیل حروف بزرگ انگلیسی به حروف کوچک انگلیسی است.

Contraction ۳.۴

یک فایل که توسط نویسنده ساخته شده (موجود در گیت‌هاب)، دارای فرمت json است. این فایل محتوی اکثر کلمات اختصار در زبان انگلیسی است. تصویر ۴ نمونه‌ای از این فایل را نشان می‌دهد. قالب کلی فایل به شکل دیکشنری هست که key همان کلمه مخفف و value شکل کامل همان کلمه مخفف است. در این قسمت از پیش‌پردازش، کلمات مختصر موجود در مجموعه داده را با شکل کامل جایگزین می‌کنیم.

English Contractions.json

```
1: {"ain't": "are not",
2: "aren't": "are not",
3: "can't've": "cannot have",
4: "can't": "cannot",
5: "'cause": "because",
6: "could've": "could have",
7: "couldn't've": "could not have",
8: "couldn't": "could not",
9: "didn't": "did not",
10: "doesn't": "does not",
11: "don't": "do not",
12: "dunno": "do not know",
13: ...}
```

شکل ۴: بخشی از کلمات اختصار زبان انگلیسی به همراه شکل کامل هر کدام

URL ۴.۴

کار بعدی، استخراج مکانیاب منبع یکسان^۱ و جابجایی آنها با یک کلید منحصر به فردی هست که بعداً در صورت نیاز می‌توان با استفاده از کلید، نشانی وب مورد نظر را بازیابی کرد. کلیدی که برای این بخش در نظر گرفته شده است دارای چنین ساختاری است:

$$\triangleright L_{counter} \triangleleft$$

به عنوان مثال، اولین نشانی وب با $\triangleright L1 \triangleleft$ جایگزین می‌شود، دومین نشانی وب با $\triangleright L2 \triangleleft$ جایگزین می‌شود و به همین ترتیب الی آخر. حرف L برگرفته از کلمه Link است. بعد از عملیات جایگزینی، نشانی وب به همراه کلید مربوطه در یک فایل ذخیره می‌شود. با این کار می‌توان نشانی وب را در صورت نیاز بازیابی کرد.

^۱Uniform Resource Locator (URL)

Email Address ۵.۴

سپس نوبت، استخراج آدرس پست الکترونیکی^۱ و جابجایی آن با یک کلید منحصر به فردی هست که بعداً در صورت نیاز می‌توان با استفاده از کلید، آدرس ایمیل موردنظر را بازیابی کرد. برای دامنه آدرس ایمیل هیچ‌گونه محدودیتی در نظر گرفته نشده است پس هر آدرس ایمیل با نام کاربری و دامنه معتبر را شناسایی می‌کنیم. کلیدی که برای این بخش در نظر گرفته شده است دارای چنین ساختاری است:

▷ $E_{counter}$ ◜

به عنوان مثال، اولین آدرس ایمیل با ▷ $E1$ ▷ جایگزین می‌شود، دومین آدرس ایمیل با ▷ $E2$ ▷ جایگزین می‌شود و به همین ترتیب الى آخر. حرف E برگرفته از کلمه Email است. بعد از عملیات جایگزینی، آدرس ایمیل به همراه کلید مربوطه در یک فایل ذخیره می‌شود. با این کار می‌توان آن را در صورت نیاز بازیابی کرد.

Currency ۶.۴

یکی دیگر از مراحل پیش‌پردازش، استخراج پول و واحد پول^۲ و جابجایی آن‌ها با یک کلید منحصر به فردی هست که بعداً در صورت نیاز می‌توان با استفاده از کلید، پول موردنظر را بازیابی کرد. کلیدی که برای این بخش در نظر گرفته شده است دارای چنین ساختاری است:

▷ $M_{counter}$ ◜

به عنوان مثال، اولین واحد پول با ▷ $M1$ ▷ جایگزین می‌شود، دومین واحد پول با ▷ $M2$ ▷ جایگزین می‌شود و به همین ترتیب الى آخر. حرف M برگرفته از کلمه Money است. بعد از عملیات جایگزینی، واحد پول به همراه کلید مربوطه در یک فایل ذخیره می‌شود. با این کار می‌توان آن را در صورت نیاز بازیابی کرد. نمونه‌هایی از اعداد و کلمه‌هایی که به عنوان پول شناسایی و با کلید مشخصی جایگزین می‌شوند، در زیر ذکر شده‌اند.

- علامت اختصاری‌هایی که هنگام استخراج در نظر گرفته می‌شوند، در جدول ۵ ذکر شده‌اند. در ضمن چه علامت قبل از عدد ظاهر شود و چه بعد از عدد، هر دو حالت از متن استخراج می‌شود.

¹Email

²Currency

#	Currency Symbol	Currency	Country
1	\$	Dollar	USA
2	€	Euro	Euro Member Countries
3	£	Pound Sterling	UK
4	¥	Yen	Japan
5	₹	Rupee	India
6	¢	Cent	USA

جدول ۵: علائم اختصاری واحد پول کشورها

- معمولا هنگام نوشتن پول به صورت عددی، اعداد سه رقم سه رقم جدا می‌شوند. اینجا کاراکتر جدا کننده می‌تواند نقطه (.) یا کاما (,) یا (') single quotation باشد.

- معمولا هنگام بیان کردن پول، عبارت‌هایی مثل «هزار»، «میلیون» و «میلیارد» به کار برده می‌شود. در جدول ۶، سیستم واحد شمارش را تعیین کردیم. به عنوان مثال، اگر در متن همراه با عدد حرف m یا کلمه Mill ظاهر شود (تفاوتبندی ندارد که سمت راست یا چپ عدد ظاهر شود) آنگاه آن حرف یا کلمه را میلیون در نظر می‌گیریم و عدد را به همراه آن حرف یا کلمه با کلید خاصی که در نظر گرفته شده جایگزین می‌کنیم.

با در کنار هم قرار دادن سه مورد قبلی، انواع حالات‌ها برای استخراج پول از متن بدست می‌آید. در ادامه چندین مثال را مشاهده می‌کنید:

€200 \$800.00 \$1,222 8\$ \$60k £300,75 ¢400'85 €600.90m €1000.10mill
2000.20Million€ \$1800'12b

#	Unit Symbol	Unit
1	K k	Thousand
2	M, Mill, Million m, mill, million	Million
3	B, Bill, Billion b, bill, billion	Billion
4	T, Trill, Trillion t, trill, trillion	Trillion

جدول ۶: سیستم شمارش پول

Time ۷.۴

استخراج و جایگزینی زمان (ساعت) یکی دیگر از مراحل پیش‌پردازش است. کلید منحصر به‌فردی که برای این قسمت در نظر گرفته شده است داری چنین ساختاری است:

$\triangleright T_{counter} \triangleleft$

به عنوان مثال، اولین ساعت با $\triangleleft T1 \triangleright$ جایگزین می‌شود، دومین ساعت با $\triangleleft T2 \triangleright$ جایگزین می‌شود و به همین ترتیب الی آخر. حرف T برگرفته از کلمه Time است. بعد از عملیات جایگزینی، زمان به همراه کلید مربوطه در یک فایل ذخیره می‌شود. با این کار می‌توان آن را در صورت نیاز بازیابی کرد. نمونه‌هایی از ساعت که به عنوان زمان شناسایی و با کلید مشخصی جایگزین می‌شوند، در جدول ۷ ذکر شده‌اند.

Type	Numbers	Postfix
Short Time	Any number separated by dot(.)	AM PM am pm
Full Time	Numbers [0:23] : Numbers [0:59] : Numbers [0:59]	A.M. P.M. a.m. p.m.

جدول ۷: سیستم تشخیص ساعت

Repeated Letter ۸.۴

معمولًا هنگام چت بعضی از حروف را به صورت تکراری می‌نویسند که می‌تواند دلایل متفاوتی داشته باشد، مثل:

- Michael Erard : «جبران کمبود نشانه‌های صوتی هنگام نوشتن به جای صحبت کردن»
- ابراز کنایه یا احساسات منفی
- مقاعده‌سازی: از نظر علمی ثابت شده است که تکرار اطلاعات احتمال تغییر ذهن افراد را افزایش می‌دهد.

یکی دیگر از مراحل پیش‌پردازش یافتن و حذف این حروف اضافی از متن هست. در زبان انگلیسی کلماتی داریم که دو حرف پشت سر هم تکرار شدند و آن کلمه معنادار است (مثلا Good) به همین دلیل مبنای تعداد حروف پشت سر هم، حداکثر دو است. پس اگر حرفی بیش از دو بار متوالی تکرار شده باشد، کل آن حروف را با دو حرف تعویض می‌کنیم، به عنوان مثال:

Goooooooood → Good

دقت بفرمایید که در این قسمت، تنها حروف انگلیسی چک نمی‌شوند بلکه کاراکتر Space نیز چک می‌شوند و اگر در متن مجموعه‌داده، بیش از دو فاصله متوالی باشد، همگی با دو فاصله جایگزین می‌شوند. در نتیجه فضاهای خالی موجود در متن نیز در این قسمت فیلتر می‌شوند. به عنوان مثال:

This_is_a_____big_____window. → This_is_a__big__window.

۵ بسط دادن

بخش سوم، مربوط به گستردگی کردن یا توسعه دادن برخی کلمات و یا کاراکترها است.

Mention ۱.۵

طبق تعریف شرکت ایکس: «مِشن توییتی شامل نام کاربری حسابی دیگر است که قبیل از آن، از علامت آت ساین (@) استفاده شده است». در این قسمت از پیش‌پردازش، هدف گسترش نام کاربری است. معمولًا نام کاربری می‌تواند یک کلمه یا ترکیبی از چندین کلمه باشد که به وسیله خط فاصله^۱ از

^۱Underline

همدیگر جدا شدند. هدف در این قسمت ابتدا شناسایی نام کاربری به کمک علامت @ است. سپس چک می‌کنیم که اگر از چندین کلمه و خط فاصله استفاده شده باشد، کلمات را به وسیله خط فاصله از همدیگر جدا و جایگزین نام کاربری اولیه می‌کنیم و در نهایت علامت @ را حذف می‌کنیم.

Hey, @dark_web This sentence is just for test.
Hey, dark web This sentence is just for test.

Hashtag ۲۰۵

هشتگ معمولاً از یک یا چندین کلمه تشکیل شده که به وسیله خط فاصله یک کلمه واحد را تشکیل می‌دهند. علامت مخصوص هشتگ «#» است. در این بخش از پیش‌پردازش، هدف پیدا کردن هشتگ‌ها و جدا کردن کلمات (در صورت امکان) و سپس جایگزین کردن با کلمه اولیه و حذف علامت هشتگ است.

Hey, #dark_web_2024 This sentence is just for test.
Hey, dark web 2024 This sentence is just for test.

۶ ذخیره مجموعه داده

بعد از انجام مراحل ذکر شده در قسمت‌های قبلی با ترتیبی که در جدول ۱ آمده، یک مجموعه داده تمیز و مرتب خواهیم داشت. فایل‌های پیش‌پردازش شده به صورت جداگانه در قالب یک فایل CSV ذخیره می‌شوند، تا در مراحل بعدی کار از آن‌ها استفاده شود.

٧ مراجـع

- [1] X (formerly Twitter). Glossary – X help center. <https://help.x.com/en/resources/glossary>, 2025. Accessed: 2025-04-05.
- [2] S. J. Whitmore. tweet-collector – GitHub repository. <https://github.com/sjwhitmore/tweet-collector>, 2025.