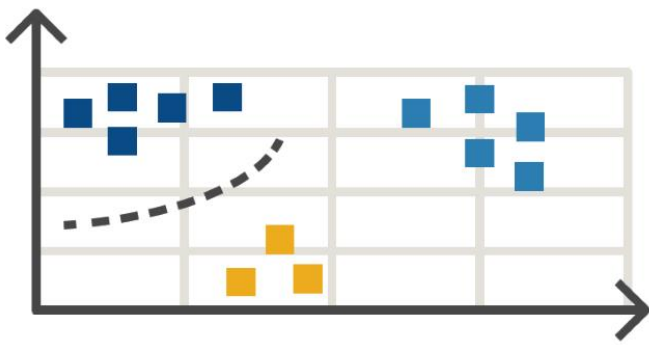


این تابع، نمودار پراکندگی دو بُعدی (2D) را از داده‌های متنی تولید می‌کند.	textscatter
این تابع، نمودار پراکندگی سه بُعدی (3D) را از داده‌های متنی تولید می‌کند.	textscatter3
این تابع، نمودار «نقشه حرارت» (Heatmap) مرتبط با عناصر داده‌های متنی را تولید می‌کند.	heatmap
این تابع، داده متنی ورودی را از طریق الگوریتم‌های استاندارد به تعدادی «دسته» (Bin) تقسیم‌بندی می‌کند. در نهایت، اطلاعاتی نظیر تعداد عناصر در هر دسته و سایر موارد، به عنوان خروجی توسط تابع تولید می‌شود.	histcounts
این تابع، داده‌های متنی (یا نمایش حاصل شده از آن‌ها) را در دسته‌ها یا طبقه‌های خاص گروه‌بندی می‌کند.	discretize

مدل‌سازی داده‌های متنی و پیش‌بینی در متلب



با استفاده از روش‌هایی نظیر Bag-of-Words و یا مدل‌های «تعبیه‌سازی کلمات» (Word Embedding) از پیش آموزش دیده، می‌توان داده‌های متنی را به «نمایش عددی» (Numerical Representation) متناظر آن‌ها تبدیل کرد.

همچنین، از طریق توابع تعریف شده در تولباکس تحلیل متن، می‌توان الگوریتم‌های یادگیری ماشین را برای پیش‌بینی و مدل‌سازی موضوعی مورد استفاده قرار داد.

توابع لازم برای مدل‌سازی داده‌های متنی و پیش‌بینی در متلب

توصیف تابع	نام تابع
با استفاده از این تابع، یک مدل «تعبیه‌سازی کلمات» (Word Embedding) از پیش آموزش دیده (که در فایل txt یا zip ذخیره شده است) توسط برنامه خوانده می‌شود.	readWordEmbedding
با استفاده از این تابع، یک مدل «تعبیه‌سازی کلمات» (Word Embedding) آموزش داده می‌شود.	trainWordEmbedding
با استفاده از این تابع، کلمات موجود در داده‌های متنی به «بردارهای تعبیه‌سازی» (Embedding Vectors) نگاشت می‌شوند.	word2vec/vec2word
مدل «تخصیص دیریکله نهان» (Latent Dirichlet allocation) یا LDA	ldaModel
مدل «تحلیل معنای نهان» (Latent Semantic Analysis) یا LSA	lsaModel
مدل Bag-of-Words یا BoW	bagOfWords
از این تابع، برای «برازش» (Fitting) مدل «تخصیص دیریکله نهان» (Latent Dirichlet allocation) یا LDA روی داده‌های متنی یا مدل نمایشی متناظر استفاده می‌شود.	fitlda

از این تابع، برای «برازش» (Fitting) مدل «تحلیل معنای نهان» (Latent Semantic Analysis) یا LSA روی داده‌های متنی یا مدل نمایشی متناظر استفاده می‌شود.	fitlsa
از این مدل، جهت پیش‌بینی موضوعات برتر موجود در اسناد متنی با استفاده از مدل تخصیص دیریکله نهان یا LDA استفاده می‌شود.	predict
از این تابع، جهت برازش توزیع احتمالی روی داده‌های متنی استفاده می‌شود.	fitdist
از این تابع، برای برازش یک مدل «رگرسیون خطی» (Linear Regression) روی داده‌های با ابعاد بالا استفاده می‌شود.	fitrlinear
از این تابع، برای برازش یک مدل «دسته‌بندی خطی» (Linear Classification) روی داده‌های با ابعاد بالا استفاده می‌شود.	fitcllinear
از این تابع، جهت برازش مدل‌های «چندکلاسه» (Multi-Class) برای دسته‌بندها (Classifiers) استفاده می‌شود.	fitcecoc

خواندن اسناد و استخراج داده‌های متنی



با استفاده از توابع تعریف شده در تولباکس تحلیل متن، این امکان برای کاربران و برنامه‌نویسان فراهم شده است تا داده‌های متنی را از فایل‌های PDF، فایل‌های Microsoft Word، فایل‌های متنی ساده و فایل‌های «صفحه گسترده» (Spreadsheet) استخراج کنند.

توابع لازم برای خواندن اسناد و استخراج داده‌های متنی

توصیف تابع	نام تابع
با استفاده از این تابع، داده‌های متنی از فایل‌های PDF، فایل‌های Microsoft Word و فایل‌های متنی ساده خوانده می‌شوند.	extractFileText
با استفاده از این تابع، «داده‌های متنی قالب‌بندی شده» (Formatted Text Data) از فایل‌ها یا رشته‌های متنی خوانده می‌شوند.	textscan
با استفاده از این تابع، ابتدا یک جدول خالی ساخته می‌شود. سپس، داده‌های متنی از فایل‌های «ستون محور» (Column-Oriented) خوانده می‌شوند. در نهایت، محتویات خوانده شده، در جدول تولید شده ذخیره می‌شوند.	readtable
از این تابع، برای تبدیل داده‌های متنی به آرایه‌های قالب‌بندی شده از نوع «رشته» (String) استفاده می‌شود.	compose



Microsoft Excel گسترده صفحه فایل‌های موجود در فایل‌های استفاده می‌شود.	xlsread
از این تابع، برای خواندن محتویات از سرویس‌های تحت وب RESTful استفاده می‌شود.	webread
از این تابع، برای ساختن یک datastore استفاده می‌شود. در متلب، یک datastore مخزنی متشکل از مجموعه‌های داده‌ای بسیار بزرگ است که نمی‌توانند در حافظه جای بگیرند. یک datastore، به کاربران و برنامه‌نویسان اجازه می‌دهد تا تمامی داده‌های خوانده و پردازش شده را در قالب یک موجودیت واحد ذخیره‌سازی و گردآوری کنند.	FileDatastore
از این تابع، جهت ساختن datastore برای فایل‌های متنی جدولی (Tabular Text Files) استفاده می‌شود.	TabularTextDatastore
از این تابع، جهت ساختن datastore برای فایل‌های صفحه گسترده (Spreadsheet Files) استفاده می‌شود.	SpreadsheetDatastore

پیش‌پردازش داده‌های متنی

“Performed preventive maintenance servicing on a broken pump.”

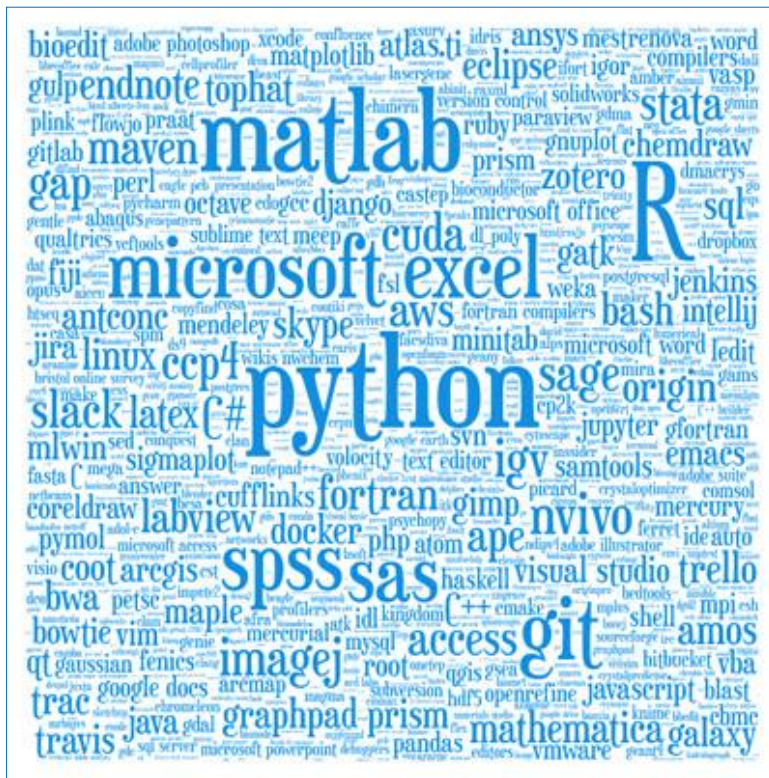
توابع تعریف شده در تولباکس تحلیل متن متلب، این امکان را برای کاربران این زبان برنامه‌نویسی فراهم می‌آورند تا «ویژگی‌هایی» (Features) را که بار معنایی ندارند و برای تحلیل (کمی و کیفی) داده‌های متنی مفید نیستند، شناسایی و حذف کنند.

ویژگی‌هایی نظیر کلمات شایع، علائم نقطه‌گذاری و آدرس‌های اینترنتی از جمله چنین ویژگی‌هایی هستند. همچنین، تکنیک‌هایی برای «نرمال‌سازی داده‌های متنی» (Text Normalization) در تولباکس تحلیل متن تعریف شده‌اند که برای بازگرداندن کلمات به «فرم ریشه‌ای» (Root Form) مورد استفاده قرار می‌گیرند (فرایند Stemming).

توابع لازم برای پیش‌پردازش داده‌های متنی

توصیف تابع	نام تابع
از این تابع، برای تقسیم‌بندی داده‌ها و اسناد متنی به مجموعه‌ای از کلمات استفاده می‌شود.	tokenizedDocument
در این تابع، از الگوریتم شناخته شده Porter جهت انجام عملیات Stemming و برگرداندن کلمات به شکل ریشه‌ای استفاده می‌شود.	normalizeWords
از این تابع، برای مدل‌سازی Bag-of-Words استفاده می‌شود.	bagOfWords

<p>در این لیست، مجموعه «کلمات بی‌اثر» (Stopwords) در زبان‌های مختلف سازمان‌دهی شده‌اند. کلمات بی‌اثر در زبان‌های مختلف، کلماتی هستند که بار معنایی ندارند و نقشی در تعیین قالب محتوایی و زمینه موضوعی اسناد ایفا نمی‌کنند.</p>	<p>stopWords</p>
<p>از این تابع، برای جستجوی ظاهر شدن یک کلمه خاص در اسناد متنی استفاده می‌شود. با استفاده از این تابع، علاوه بر آمار مرتبط با ظاهر شدن یک کلمه در اسناد متنی، «زمینه موضوعی» (Context) ظاهر شدن کلمه نیز استخراج می‌شود.</p>	<p>context</p>
<p>با استفاده از این تابع، مجموعه‌ای از کلمات انتخابی (نظیر کلمات بی‌اثر) از اسناد متنی یا مدل‌های Bag-of-Words حذف می‌شوند.</p>	<p>removeWords</p>
<p>با استفاده از این تابع، کلمات بلند (با تعداد کاراکترهای زیاد) از اسناد متنی یا مدل‌های Bag-of-Words حذف می‌شوند. پارامتر طول (تعداد کاراکترها) کلمات توسط کاربر مشخص می‌شود.</p>	<p>removeLongWords</p>
<p>با استفاده از این تابع، کلمات کوتاه (با تعداد کاراکترهای کم) از اسناد متنی یا مدل‌های Bag-of-Words حذف می‌شوند. پارامتر طول کلمات (تعداد کاراکترها) توسط کاربر مشخص می‌شود.</p>	<p>removeShortWords</p>
<p>با استفاده از این تابع، کلماتی که تعداد دفعات تکرار آن‌ها در اسناد یا داده‌های متنی یا مدل Bag-of-Words کم است، حذف می‌شوند. پارامتر تعداد دفعات تکرار مطلوب، توسط کاربر مشخص می‌شود.</p>	<p>removeInfrequentWords</p>
<p>با استفاده از این تابع، علائم نقطه‌گذاری از اسناد و داده‌های متنی حذف می‌شوند.</p>	<p>erasePunctuation</p>



نوع داده «رشته» (String)

"Hello,world"

توابع تعریف شده در تولباکس تحلیل متن، این امکان را برای کاربران و برنامه‌نویسان فراهم می‌آورند تا داده‌های متنی را به شکل بهینه‌ای «دستکاری» (Manipulate)، «مقایسه» (Compare) و «ذخیره» (Store) کنند.

توابع لازم برای دستکاری، مقایسه و ذخیره متغیرهای نوع داده «رشته» (String)

توصیف تابع	نام تابع
این دستور، چگونگی تعریف متغیر رشته‌ای (یک متغیر از نوع داده رشته) را نمایش می‌دهد.	<code>str = "Hello,world"</code>
این دستور، چگونگی تعریف آرایه رشته‌ای (یک آرایه از نوع داده رشته) را نمایش می‌دهد.	<code>str = ["Hello", "World"]</code>
از این تابع، برای تبدیل یک بردار از کاراکترها به نام C، به یک متغیر رشته‌ای استفاده می‌شود.	<code>str = string(C)</code>
از این تابع، برای تبدیل یک متغیر رشته‌ای به اعداد اعشاری (از نوع double) استفاده می‌شود.	<code>str2double</code>
با استفاده از این تابع، طول رشته‌ها در خروجی نمایش داده می‌شود.	<code>strlen</code>
این تابع، آرگومان ورودی را مورد بررسی قرار می‌دهد تا مشخص کند که این ورودی، آرایه‌ای از نوع رشته است یا نه.	<code>isstring</code>
از این تابع، برای ترکیب رشته‌ها استفاده می‌شود.	<code>join</code>
از این تابع، برای تقسیم‌بندی رشته‌ها به آرایه رشته‌ای استفاده می‌شود.	<code>split</code>
در این تابع، جهت تقسیم‌بندی رشته‌ها به آرایه رشته‌ای، باید کاراکترهای <code>newline</code> در رشته مشاهده شوند.	<code>splitlines</code>
از این تابع، برای پیدا و جا به جا کردن «زیر رشته‌ها» (Substring) در یک آرایه رشته‌ای استفاده می‌شود.	<code>replace</code>
با استفاده از این تابع، رشته‌ای که به عنوان آرگومان ورودی به تابع پاس داده می‌شود مورد بررسی قرار می‌گیرد تا مشخص شود یک زیر رشته یا «الگوی» (Pattern) خاص در این رشته وجود دارد یا نه.	<code>contains</code>
از این تابع، برای حذف کردن زیر رشته‌های درون متغیرهای رشته‌ای استفاده می‌شود.	<code>erase</code>
از این تابع، برای استخراج زیر رشته‌های موجود میان «شاخص‌های» (Indicators) تعریف شده استفاده می‌شود.	<code>extractBetween</code>



از این تابع، برای استخراج زیر رشته از یک مکان خاص به بعد، در یک متغیر رشته‌ای استفاده می‌شود.	extractAfter
از این تابع، برای استخراج زیر رشته تا یک مکان خاص، در یک متغیر رشته‌ای استفاده می‌شود.	extractBefore
از این تابع، برای مقایسه محتویات متغیرهای رشته‌ای استفاده می‌شود.	strcmp
از این تابع، برای انجام عملیات مرتبط با «عبارات منظم» (Regular Expressions) استفاده می‌شود؛ به طور ویژه، برای تطبیق دادن عبارات منظم با رشته‌ها و زیر رشته‌های موجود در آن‌ها. شایان ذکر است که عملیات عبارات منظم روی رشته‌ها در زبان متلب، به کوچک یا بزرگ بودن حروف حساس است (Case Sensitive).	regexp

مجموعه آموزش‌های داده‌کاوی و یادگیری ماشین (+کلیک کنید)

برای مشاهده دیگر «تقلب‌نامه‌های» مجله فرادرس، به [این لینک](#) مراجعه فرمایید.

جهت آگاهی از آخرین تقلب‌نامه‌های منتشر شده، در [کانال تلگرام](#) مجله فرادرس عضو شوید.

تهیه و تنظیم: مجله فرادرس

