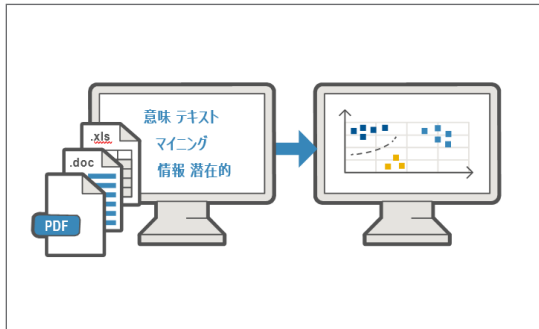


Text Analytics Toolbox のクイックスタートガイド



Text Analytics Toolbox™ は、テキストデータの前処理、解析、モデリングのためのアルゴリズムと可視化機能を提供します。ツールボックスを使用して作成されたモデルは、感情分析、予知保全、トピックモデリングなどのアプリケーションで使用できます。

関連情報: mathworks.co.jp/products/text-analytics

関数名	説明
<code>wordcloud</code>	bag-of-words あるいは LDA モデルからのワードクラウド作成
<code>wordCloudCounts</code>	ワードクラウド作成のための単語数カウント
<code>textscatter</code>	テキストの散布図
<code>textscatter3</code>	テキストの3次元散布図
<code>heatmap</code>	ヒートマップチャートの作成
<code>histcounts</code>	ヒストグラムのピンのカウント数
<code>discretize</code>	データをピンまたはカテゴリにグループ化

可視化

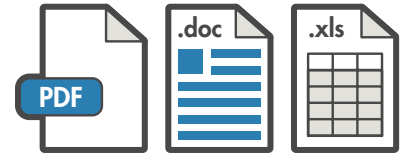
ワードクラウドとテキスト散布図を使用して、結果を要約して検証します。

モデリングと予測

Bag-of-words や学習済みの単語分散表現モデルを使用してテキストを数値表現に変換し、予測やトピックモデリングに特化した機械学習アルゴリズムを適用します。

関数名	説明
<code>readWordEmbedding</code>	テキストファイルからの単語の分散表現の読み込み
<code>trainWordEmbedding</code>	単語の分散表現を学習
<code>word2vec/vec2word</code>	単語を分散ベクトルにマップ
<code>ldaModel</code>	潜在的ディリクレ配分法 (LDA) モデル
<code>lsaModel</code>	潜在意味解析 (LSA) モデル
<code>bagOfWords</code>	Bag-of-words モデル
<code>fitlda</code>	潜在的ディリクレ配分法 (LDA) モデルの学習
<code>fitlsa</code>	潜在意味解析 (LSA) モデルの学習
<code>predict</code>	文書内の主要な LDA トピックを予測
<code>fitdist</code>	データへの確率分布オブジェクトの近似
<code>fitrlinear</code>	高次元データに対する線形回帰モデルのあてはめ
<code>fitclinear</code>	高次元データに対する線形分類モデルのあてはめ
<code>fitcecoc</code>	サポート ベクター マシンまたはその他の分類器向けのマルチクラスモデルの近似

関数名	説明
<code>extractFileText</code>	PDF, Microsoft Word, テキストファイルからのデータ読み込み
<code>textscan</code>	テキストファイルまたは文字列からの書式付きデータの読み込み
<code>readtable</code>	ファイルからのテーブルの作成
<code>compose</code>	データを書式設定された string 配列に変換
<code>xlsread</code>	Microsoft Excel スプレッドシート ファイルの読み取り
<code>webread</code>	RESTful Web サービスからのコンテンツの読み取り
<code>TabularTextDatastore</code>	表形式テキスト ファイルのデータ ストア
<code>FileDatastore</code>	カスタム ファイル リーダーを使用するデータ ストア
<code>SpreadsheetDatastore</code>	スプレッドシート ファイルのデータ ストア



読み込み

Microsoft® Word® ファイル、PDF、テキストファイル、スプレッドシートからテキストを読み込みます。

✦ 破損したポンプで
予知保全サービスを
実施した

前処理

ありきたりな単語、句読点、および URL などの深い意味を持たない単語を取り除き、語幹抽出により単語の正規化を行います。

関数名	説明
<code>tokenizedDocument</code>	文書を単語の集合に分割
<code>normalizeWords</code>	Porter stemmer を使用して語形変化で使用される語尾を単語から削除
<code>bagOfWords</code>	Bag-of-words モデル
<code>stopWords</code>	ストップワードリスト
<code>context</code>	文書内の単語の出現を検索
<code>removeWords</code>	文書あるいは bag-of-words から特定の単語を削除
<code>removeLongWords</code>	文書あるいは bag-of-words から長い単語を削除
<code>removeShortWords</code>	文書あるいは bag-of-words から短い単語を削除
<code>removeInfrequentWords</code>	bag-of-words モデルから頻度の低い単語を削除
<code>erasePunctuation</code>	テキストや文書から句読点の削除

関数名	説明
<code>str = "Hello,world"</code>	string 変数を定義
<code>str = ["Hello", "World"]</code>	string 配列を定義
<code>str = string(C)</code>	キャラクタ配列 C を string 型に変換
<code>str2double</code>	文字列の double 型への変換
<code>strlen</code>	文字列の長さ
<code>isstring</code>	入力が string 配列かどうかを判別
<code>join</code>	文字列の結合
<code>split</code>	文字列の分割
<code>splitlines</code>	文字列を改行文字の位置で分割
<code>replace</code>	文字列内の部分文字列を検索して置換
<code>contains</code>	文字列内に指定したパターンがあるかどうかの判別
<code>erase</code>	文字列内の部分文字列の削除
<code>extractBetween</code>	部分文字列の開始と終了を示す指定子の間の部分文字列を抽出
<code>extractAfter</code>	指定された位置より後の部分文字列の抽出
<code>extractBefore</code>	指定した位置より前の部分文字列を抽出
<code>strcmp</code>	文字列の比較
<code>regexp</code>	正規表現の一致 (大文字小文字を区別)

"Hello,world"

String 配列

テキストデータを効率的に操作、比較、保存します。
(各言語に対応)