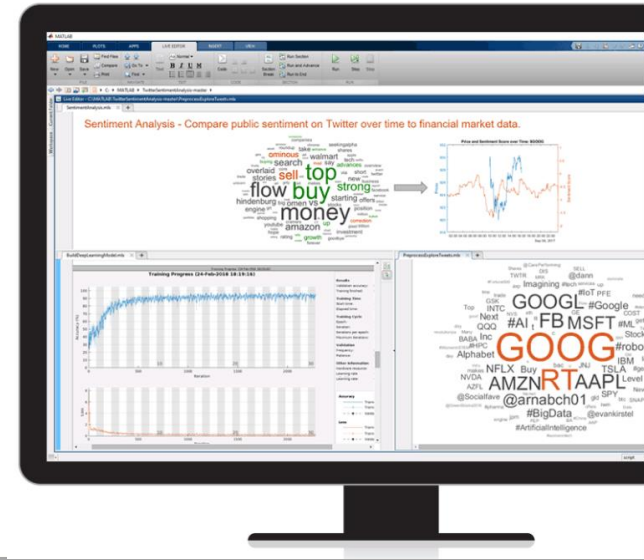# 비정형 데이터의 숨어있는 가치 창출을 위한 Text Analytics

*A hands-on MATLAB Workshop*

June. 3rd, 2020

**송완빈 과장**
Wanbin Song
Application Engineer @ MathWorks
wsong@mathworks.com

For Online Workshop

# About this workshop

- ## What this workshop is about
  - Text Analytics Toolbox overview – For Korean language
  - Hands-on introduction to the tools
    - We will not cover ***all*** the capabilities

- ## What this workshop is NOT about
  - Expertise in Text Analytics / Natural Language Processing
  - Expertise with Text Analytics Toolbox

# Way of working – Rules of the Road

- Slides to introduce topics

- Examples to demonstrate functionality

- Exercises to *get your hands dirty* and try something beyond what is demonstrated

  – Suggestion: do not modify the functions but save them with different names so that you can trace your steps back

  – This is intended for learning so…don't look at solutions unless you get behind

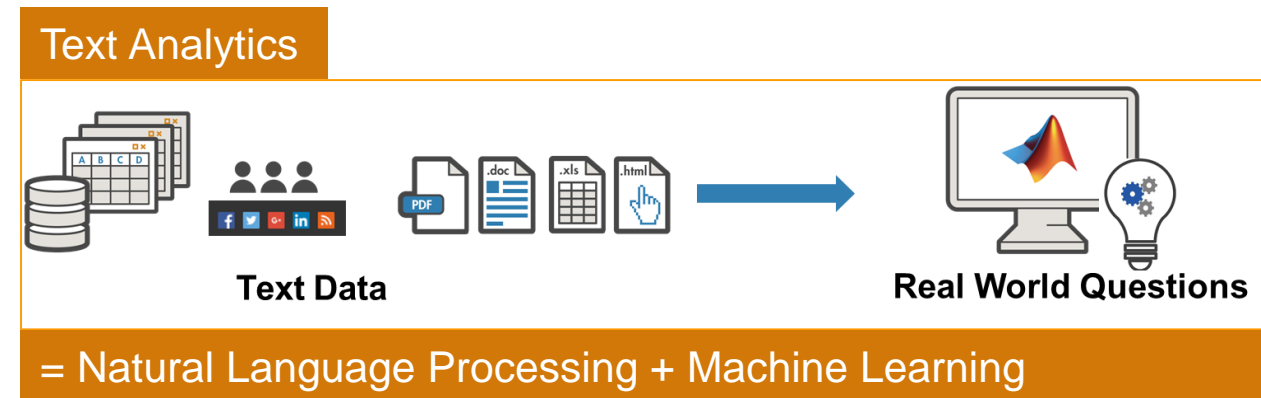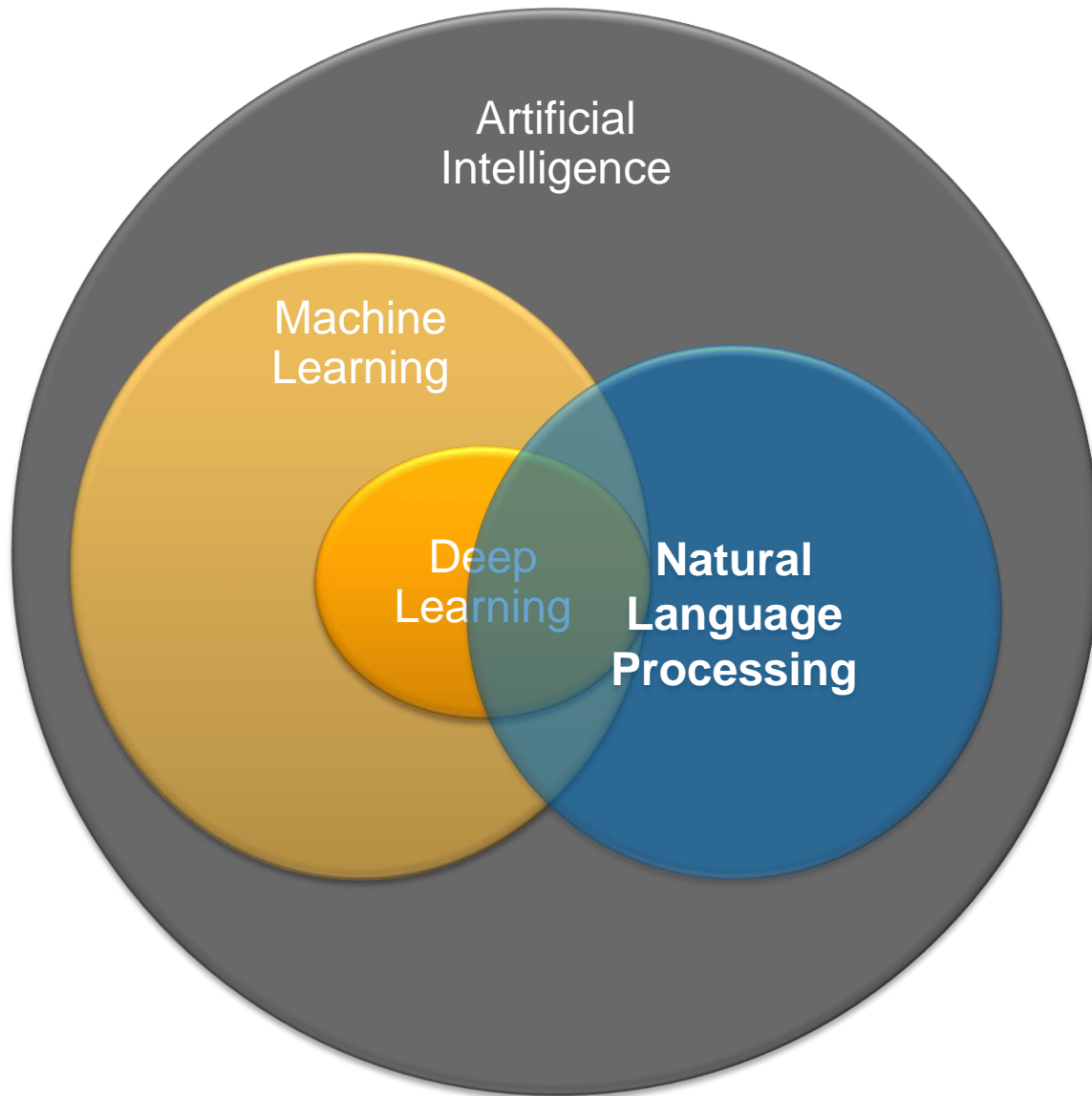  – Each exercise starts anew, so don't worry if you get behind

**>> Be curious and ask all the questions you like!**

*Raise Your Hand!( Not the real one )*  JL  Jiyoun Lee

# What is Natural Language Processing? Text Analytics?



Artificial Intelligence

Machine Learning

Deep Learning

**Natural Language Processing**

Text Analytics

Text Data

Real World Questions

= Natural Language Processing + Machine Learning

# What Makes It Difficult?

- Many words in a language, same word different meaning, dialects

- Machines understand logic.. Human beings not so logical

- Ambiguity, emotion, subjectivity, personality, culture …

연예인 고수?

쌀국수 전문가?

쌀국수는 고수랑 먹어야 제 맛이지!

향이 강한 고수?

그래 너 잘났다~ 아주 대단하셔요 전문가 납셨어~
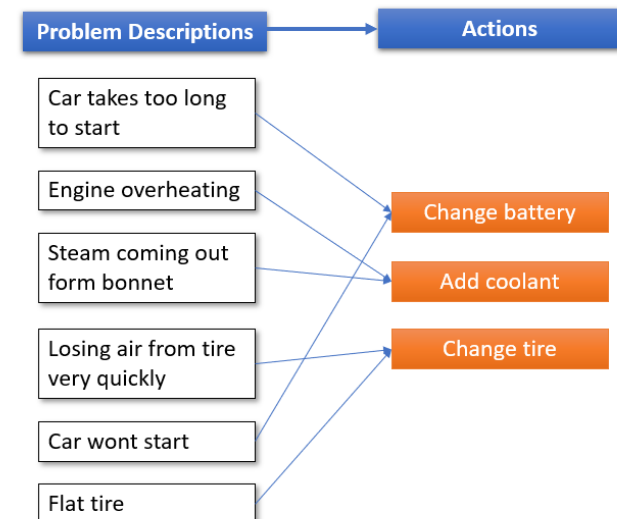
# What Is Text Analytics Used for?

- **Topic Modeling**

  *Identify topics from a collection of documents that shows underlying patterns and relationships in raw text data.*
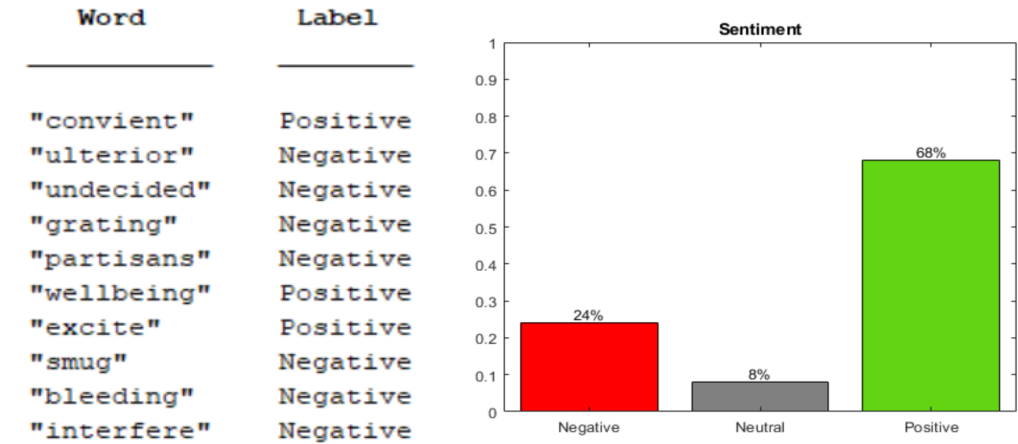
- **Text Classification**

  *Classify documents into pre-determined categories for efficient information retrieval and prediction.*
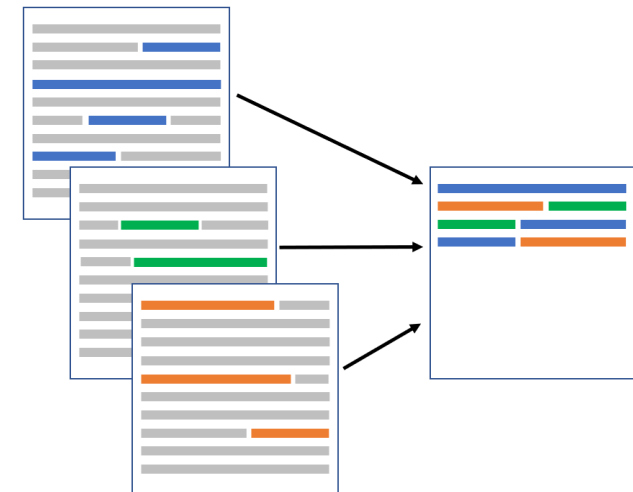
# What Is Text Analytics Used for?

▪ **Sentiment Analysis**
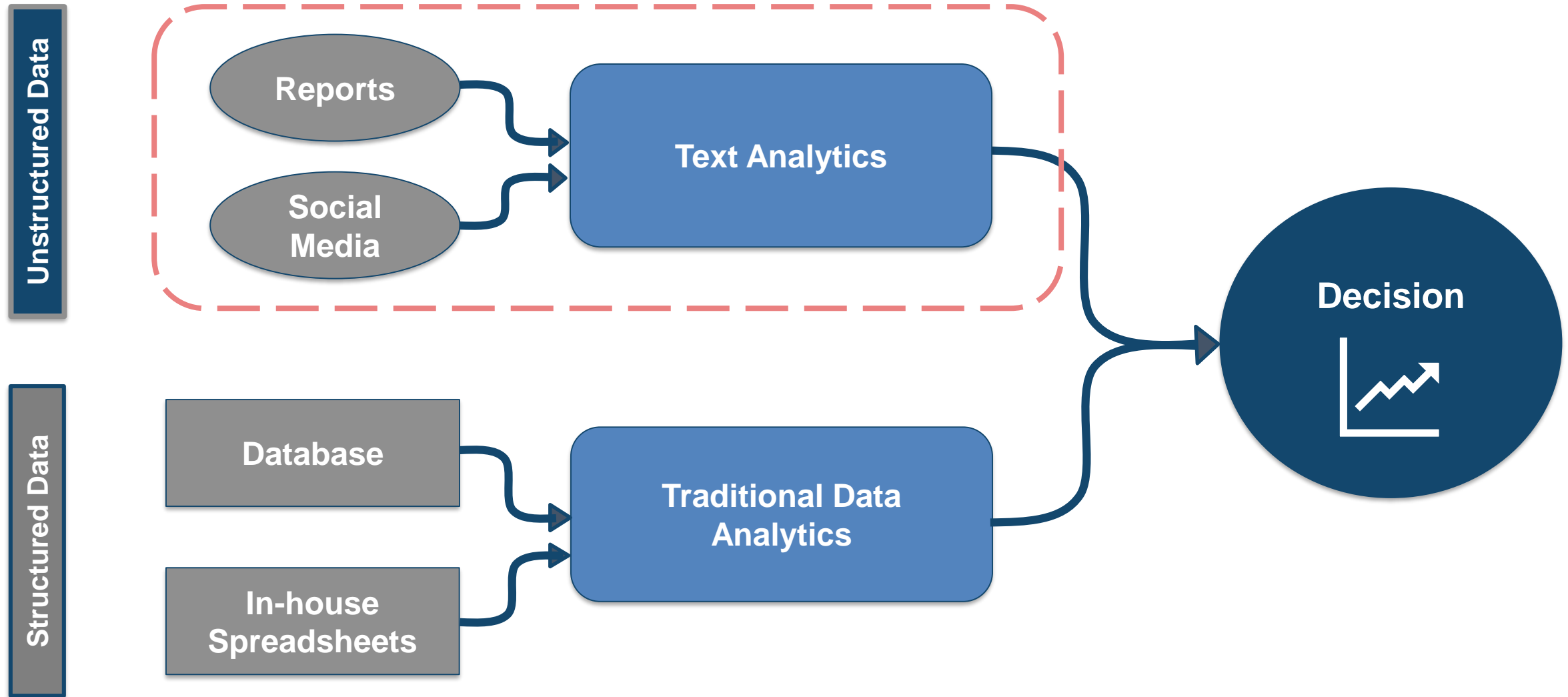
 *Identify and score sentiments expressed in text.*

| Word | Label |
| --- | --- |
| "convient" | Positive |
| "ulterior" | Negative |
| "undecided" | Negative |
| "grating" | Negative |
| "partisans" | Negative |
| "wellbeing" | Positive |
| "excite" | Positive |
| "smug" | Negative |
| "bleeding" | Negative |
| "interfere" | Negative |

▪ **Summarization**

 *Extract a summary from one or more documents automatically.*

# Big Picture

# Text Analytics Workflow

**1** Access and Explore → **2** Preprocess → **3** Develop Predictive Models → **4** Share Insights and Models

## 1 — Access and Explore

**Data Sources**

**Languages**

**Visualize**

## 2 — Preprocess

**Raw Data**

**Clean Data**

## 3 — Develop Predictive Models

**Deep Learning**

**Machine Learning**

**Statistics**

## 4 — Share Insights and Models

**Report, Publish, Interactive Notebook**

**Desktop & Web Apps**

**Production & Web App Servers**

# Strings
*The better way to work with text*

- Manipulate, compare, and store text data efficiently

```
>> "image" + (1:3) + ".png"

   1×3 string array

      "image1.png"    "image2.png"    "image3.png"
```

- Simplified text manipulation functions

```
methods string
```

string 클래스에 대한 메서드:

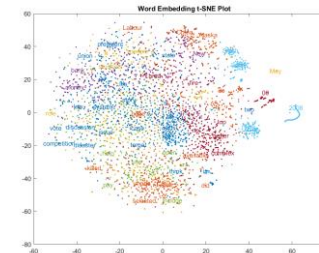| | | | | | |
|---|---|---|---|---|---|
| cellstr | eq | gt | lower | reverse | upper |
| char | erase | insertAfter | lt | sort | |
| compose | eraseBetween | insertBefore | ne | split | |
| contains | extractAfter | ismissing | pad | splitlines | |
| count | extractBefore | issorted | plus | startsWith | |
| double | extractBetween | join | replace | strip | |
| endsWith | ge | le | replaceBetween | strlength | |

# ① Access and Explore Text Data

- Import text from databases, social media, news feeds, equipment logs, reports, and surveys

- Visually explore data with word clouds & scatter plots

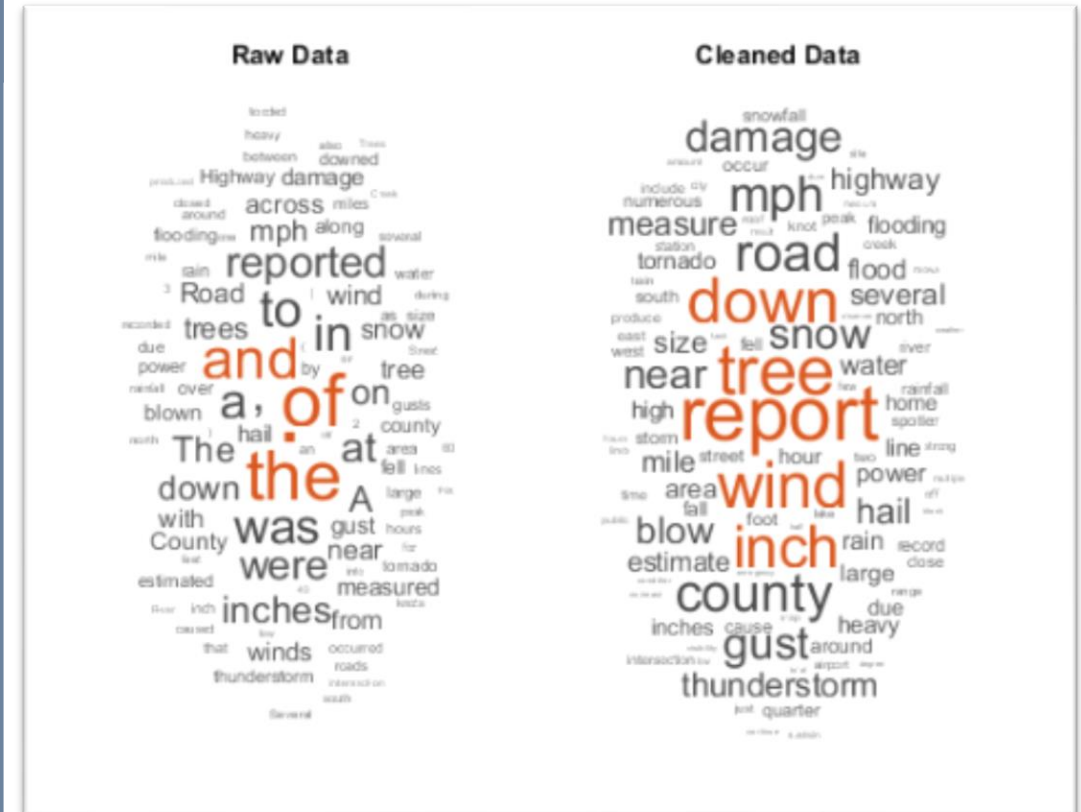- Local language support for Japanese(18b), German (19a), and Korean (19b)

# **②** Preprocess Data
## *Prepare Text Data for Model Building*

| Text data may contain: |
|---|
| ▪ Variations in case, for example "new" and "New" |
| ▪ Variations in word forms, for example "walk" and "walking" |
| ▪ Words which add noise, for example stop words such as "the" and "of" |
| ▪ Punctuation and special characters |
| ▪ HTML and XML tags |

# ③ Build Predictive Models with Text

## Convert Text to Numbers

- Bag of Words or N-grams

- Term Frequency-Inverse Document Frequency

- Word Embedding (FastText, Glove, train your own or read in someone else's)
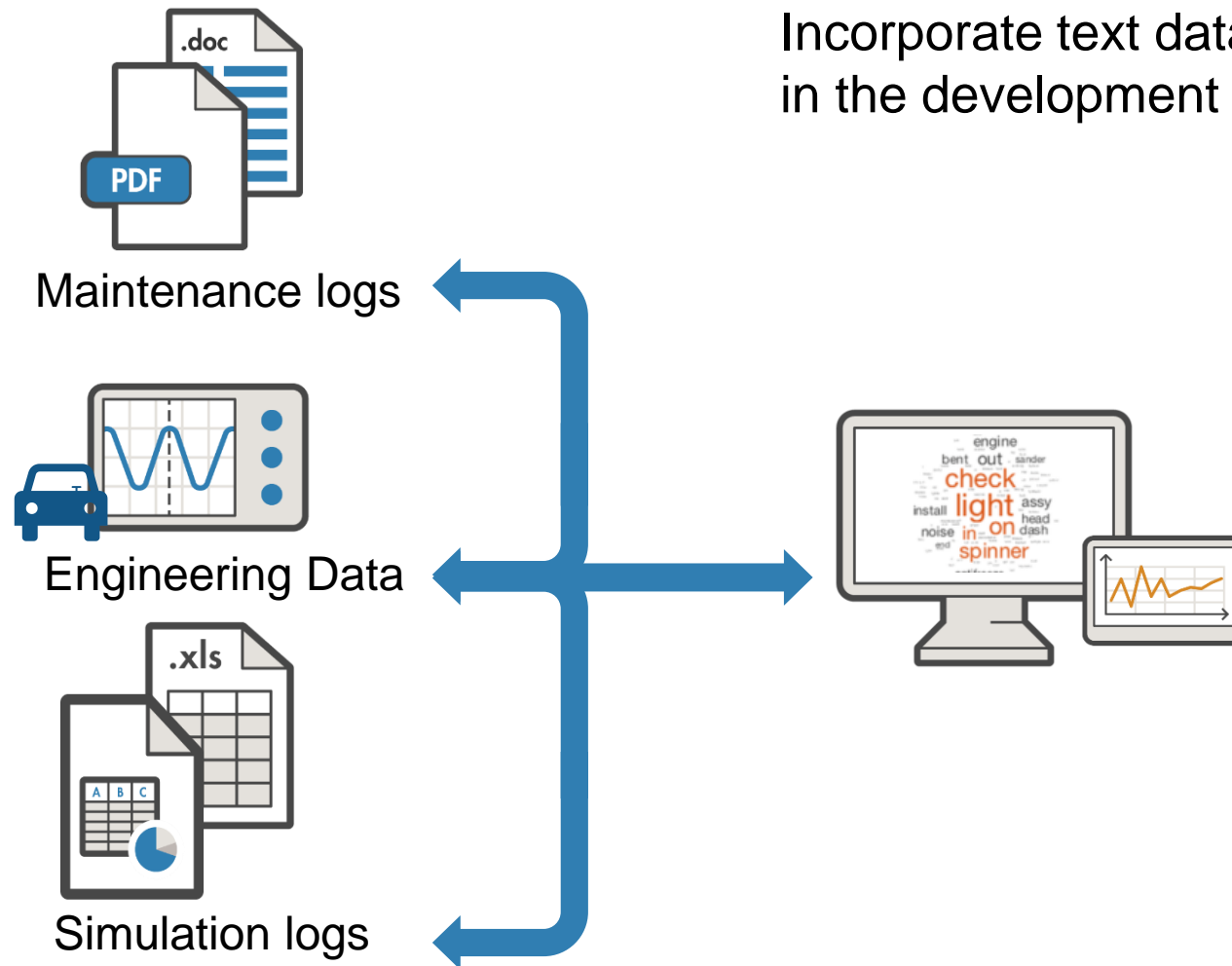
- Dimensionality Reduction

## Build Models with Machine & Deep Learning Algorithms

- Choose an algorithm based on application

- State-of-the-art machine and deep learning algorithms available in MATLAB
  - Classification
  - Clustering
  - Descriptive Statistics
  - Regression

**③ DEMO:** Topic Modeling

**Goal :** Identify key topics in documentations.

Incorporate text data with other types of engineering data in the development of smart systems.



Maintenance logs

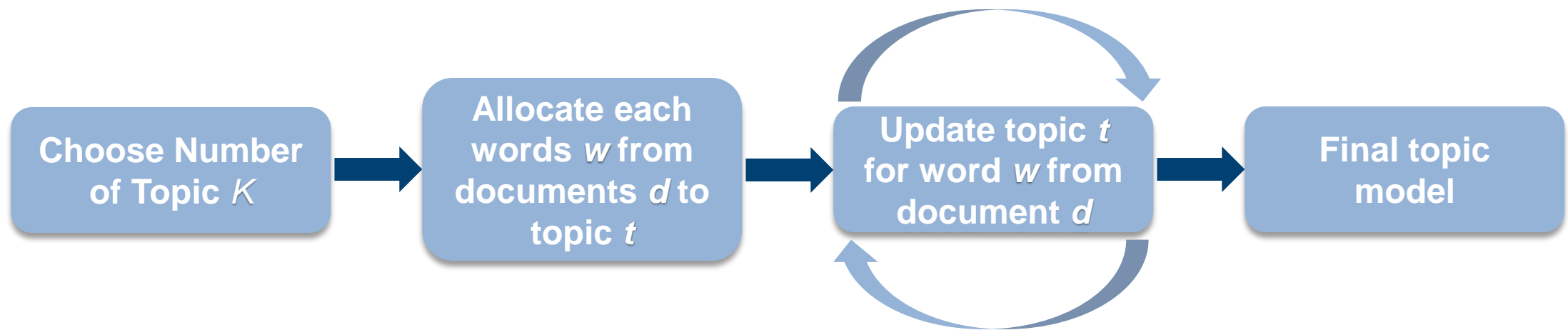Engineering Data

Simulation logs

*Sample Topics identified:*

# LDA (Latent Dirichlet Allocation)

- **Topic modeling** is a type of statistical **modeling** for discovering the abstract "**topics**" that occur in a collection of documents.

- Latent Dirichlet Allocation (**LDA**) is an example of **topic model** and is used to classify text in a document to a particular **topic**.

| Choose Number of Topic $K$ | → | Allocate each words $w$ from documents $d$ to topic $t$ | → | Update topic $t$ for word $w$ from document $d$ | → | Final topic model |
|---|---|---|---|---|---|---|

- $P(topic\ \boldsymbol{t}\ |document\ \boldsymbol{d})$
- $P(word\ \boldsymbol{w}\ |topic\ \boldsymbol{t})$

# ③ **DEMO:** Sentiment Analysis

**Goal** : Determining real-time sentiment scores for use in financial trading strategies

| Positive (+) | Neutral | Negative (-) |
|---|---|---|
| (growth, advances, up, strong) | **Hold** | (bust, difficulty, lack, struggle) |
| **Buy!** | | **Sell!** |

**Other Applications:**

- Automating the classification of reviews, whether positive or negative
- Analyzing surveys to understand why customers are satisfied or dissatisfied
- Assessing counterparty credit risk

# Demo: Workflow

**Data**

```
ans = 508×1 string array
    "Walmart: "you wanna destroy Amazon?" Google: "bet" $WMT $GOOG
    "$WMT wants next level customer service w/highly personalized
    "Ironic prelude to $DIS buying $TWTR soon IMO $AAPL $GOOG $SPY
    "$AMZN the $WMT threat grows each and every day  https://t.co/
    "MU Investments Co. Ltd. Sells 30 Shares of Alphabet Inc. $GOO
    "Ad $ are going to $GOOG and $FB away from wppgy  #Advertising
    "Big bullish unusual option activity detected: $SPX, $GOOG, $G
    "REPORT: Apple to build data center in Iowa: https://t.co/jwH6
    "RT @theflynews: REPORT: Apple to build data center in Iowa: h
```

**Preprocess text**

**Test data**

**Trained model**

**Score**

**Word Embedding**

```
wordEmbedding with properties:

    Dimension: 100
    Vocabulary: [1×1193514 string]
```
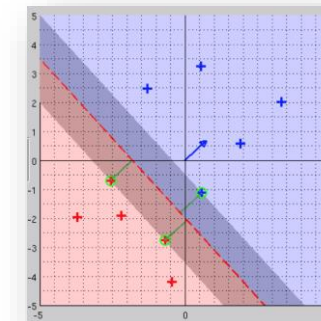
**Positive + Negative Word List**

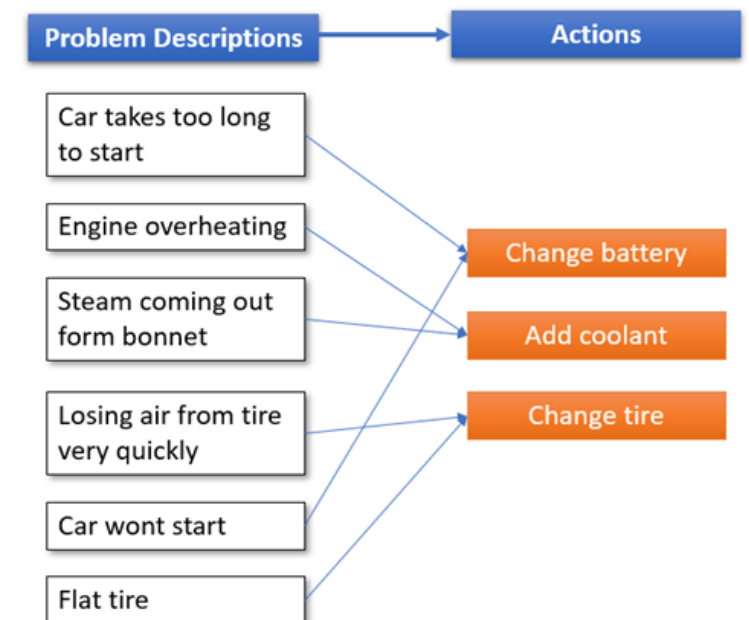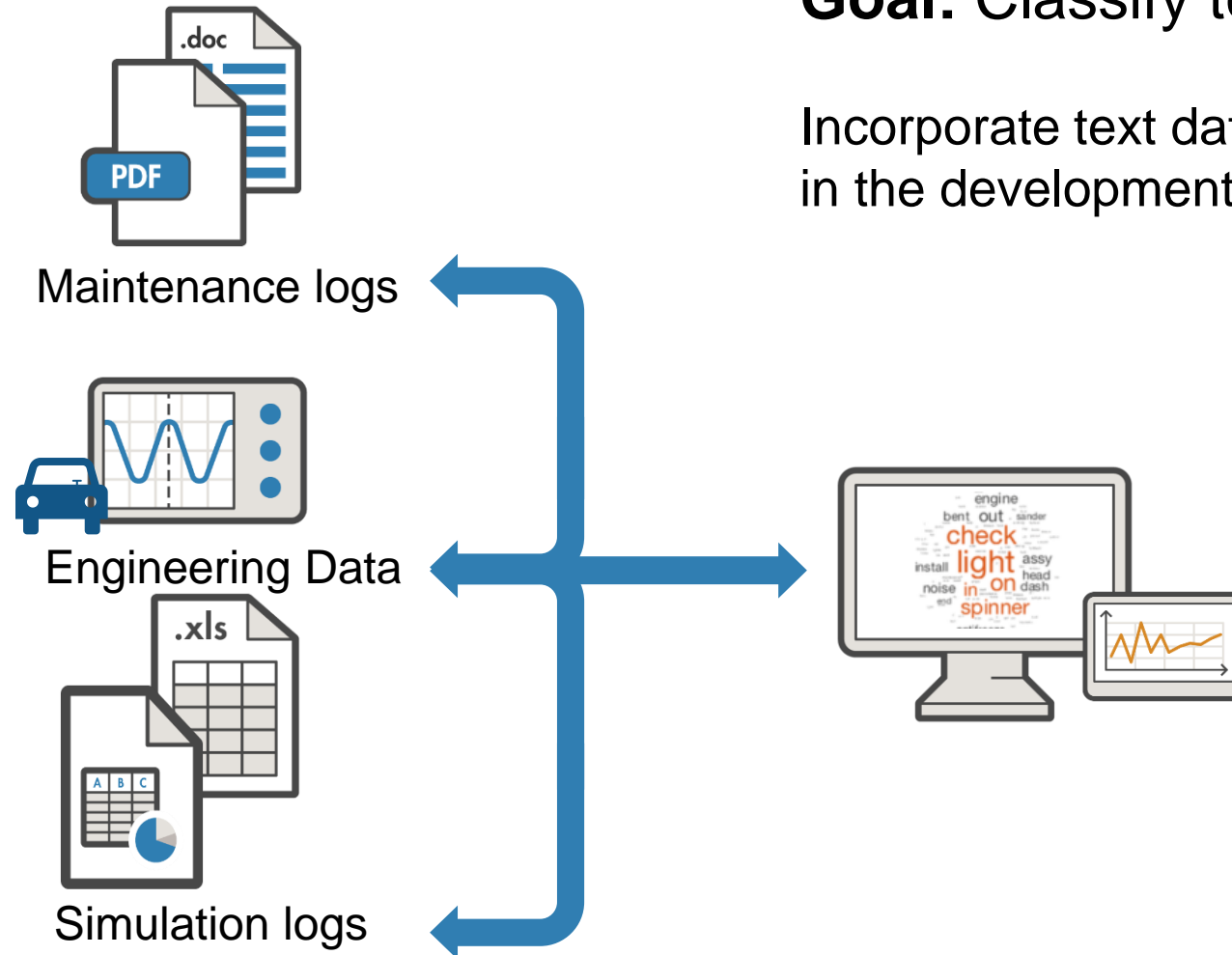| pos | | neg | |
|---|---|---|---|
| 2006x1 string | | 4783x1 string | |
| | 1 | | 1 |
| 1 | a+ | 1 | 2-faced |
| 2 | abound | 2 | 2-faces |
| 3 | abounds | 3 | abnormal |
| 4 | abundance | 4 | abolish |
| 5 | abundant | 5 | abominable |
| 6 | accessable | 6 | abominably |
| 7 | accessible | 7 | abominate |

**Training data**

**Machine Learning Model**

# ③ **DEMO:** Document Classification

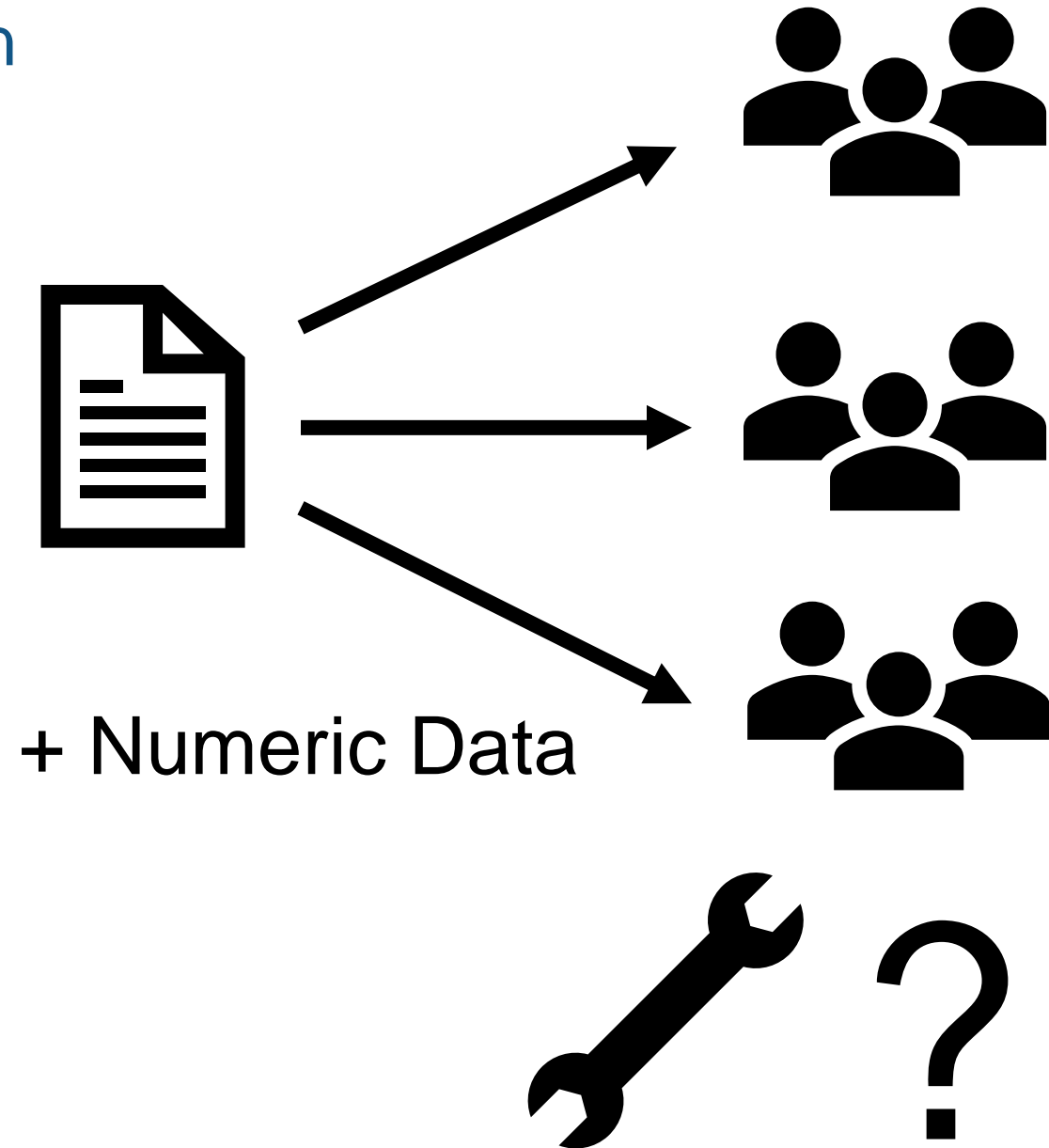**Goal:** Classify text based on historical data

Incorporate text data with other types of engineering data in the development of smart systems.

Maintenance logs

Engineering Data

Simulation logs

**Problem Descriptions** → **Actions**

Car takes too long to start

Engine overheating

Steam coming out form bonnet

Losing air from tire very quickly

Car wont start

Flat tire

Change battery

Add coolant

Change tire

**3** **DEMO:** Document Classification

- Document Classification
  - Field reports
  - Bug reports

- Predictive Maintenance
  - Equipment log notes
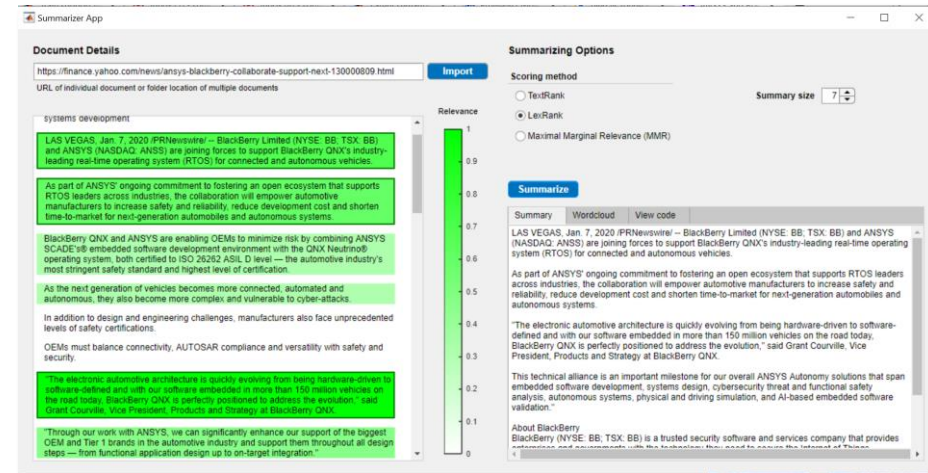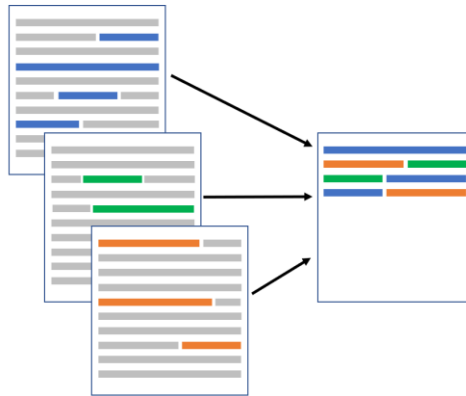
+ Numeric Data

**4** Share Insights and Models
*Many options for MATLAB and non-MATLAB users*

- Share interactive Live Editor Notebook

- Publish scripts as HTML, PDF, LaTeX, or Microsoft Word

- Generate formatted reports using MATLAB Report Generator

- Create standalone or web apps using MATLAB App Designer

- Host the application on MATLAB Production Server or MATLAB Web App Server

# Document Summarization
*Extract a summary from one or more documents automatically*

[Extractive Summarization](#)



## Relevant Applications

- Identify opportunities and gaps in scientific research by summarizing technical articles.

- Highlight and understand relevant information faster by summarizing internal reports.

# Spelling Correction

[correctSpelling](correctSpelling)

```
str = [
    "A documnent containing some misspelled worrds."
    "Another documnent cntaining typos."];
documents = tokenizedDocument(str);
```

⬇

```
updatedDocuments = correctSpelling(documents)
```

⬇

```
updatedDocuments =
  2×1 tokenizedDocument:

    7 tokens: A document containing some misspelled words .
    5 tokens: Another document containing typos .
```
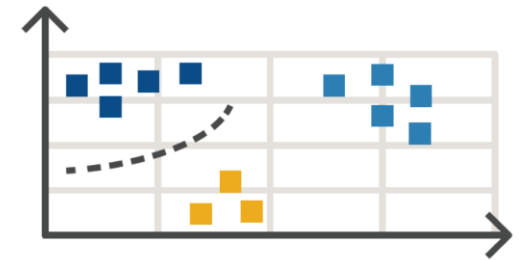
*NOTE: Supports English, German, and Korean text.*

# Text Analytics Toolbox

**Import**

**Preprocess**

**Visualize**

**Model and Predict**

- Extract text from Microsoft Word files, PDFs, text files and spreadsheets.

- Remove less helpful artifacts and apply text normalization.

- Use word clouds and text scatter plots to summarize and validate results.

- Convert text into numeric representations and apply specialized machine learning algorithms for prediction and topic modeling.

# More Resources



- [Text Analytics Toolbox](#)

- [Text Analytics Toolbox – Documentation](#)

- [Getting Started with Text Analytics in MATLAB (White Paper)](#)

- [Text Analytics in MATLAB (23:36) – Video](#)

- [8 MATLAB Cheat Sheets for Data Science - Cheat Sheets](#)

**Visit MathWorks at**