

A graphic of the MATLAB logo, consisting of five overlapping triangles in shades of blue and orange, positioned on the left side of the image.

MATLAB EXPO 2018  
KOREA

# MATLAB EXPO 2018

비정형 데이터의 숨어있는 가치  
창출을 위한 Text Analytics

송완빈

Application Engineer

MathWorks Korea



# Agenda

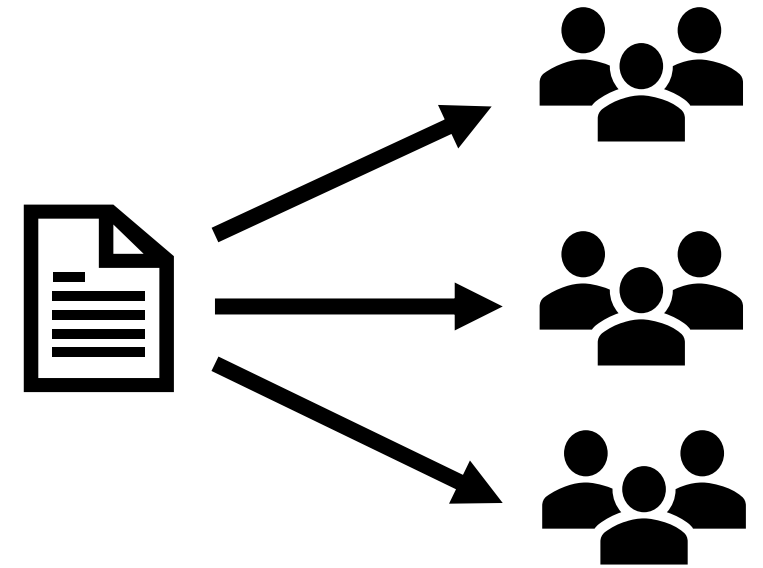
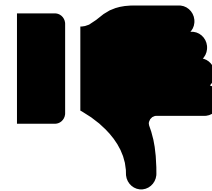
- What is Text Analytics?
- Text Analytics Workflow
- Simple Demo for Text Analytics Technique
- Sentiment Analysis with Text Analytics
- Key Takeaways

# What is Text Analytics?

- Text Analytics
  - The goal of deriving information from text data
- Text Mining
  - Older phrase for ‘Text Analytics’
- NLP(Natural Language Processing)
  - A method which leverages human language (syntax, semantics, discourse, speech)

# Text As Data

- Document Classification
  - Field reports
  - Bug reports
- Sentiment Analysis
  - Survey data
  - Trial notes
  - Social media
- Predictive Maintenance
  - Equipment log notes



+ Numeric Data

## Text Analytics Toolbox

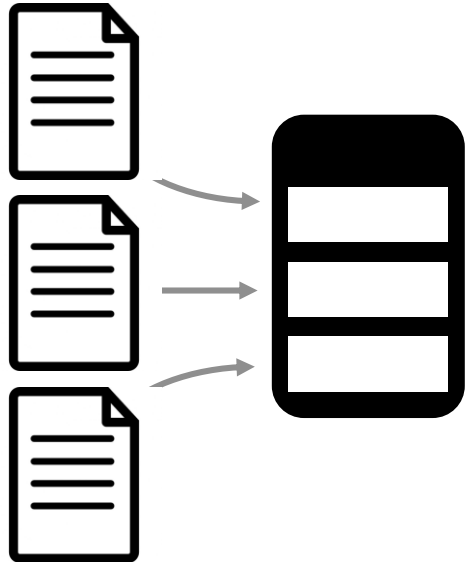
Access and Explore Data

Preprocess Data

Develop Predictive Models

Clean-up Text

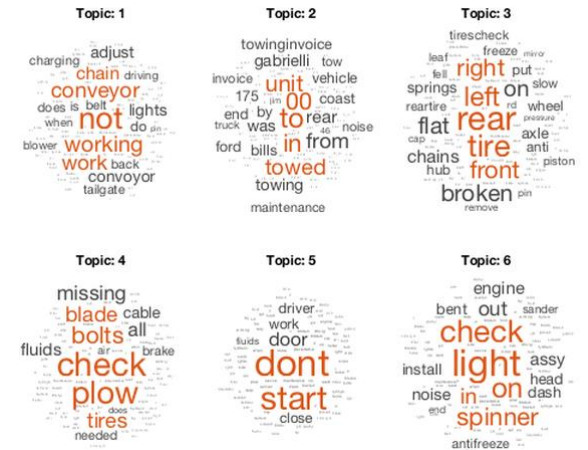
Convert to Numeric



Media reported two trees blown down along I-40 in the Old Fort area.

media report two tree blown down i40 old fort area

	cat	dog	run	two
doc1	1	0	1	0
doc2	1	1	0	1



- Word Docs
- PDF's
- Text Files
- HTML Files

- Stop Words
- Stemming
- Tokenization

- Bag of Words
- TF-IDF
- Word Counting

- Word-Embedding
- Machine Learning
  - LDA
  - LSA
- Deep Learning
  - LSTM

# Strings

## *The better way to work with text*

**R2016b**

- 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

- Example: Check if a string is contained within another string

- Previously: `if ~isempty(strfind(textdata, "Dog"))`
- Now: `if contains(textdata, "Dog")`

- Performance improvement

- Up to 50x faster using **contains** with **string** than **strfind** with **cellstr**
- Up to 2x memory savings using **string** over **cellstr**

# Brief overview on Toolbox with Simple Demo



# Text Analytics Toolbox

# R2017b

## ■ Sources of Text Data

- Maintenance Logs
- News/Social Media
- Customer Surveys
- Field Reports
- Research Papers

### Split Text into Individual Words

```
documents = tokenizedDocument(repairNotes)
```

### Create a Bag-of-Words Model

```
bag = bagOfWords(documents)
```

### Fit a Topic Model with 4 Topics

```
numTopics = 4;
mdl = fitlda(bag,numTopics)
```

```
(36,1) coolant leak spinner light outwip
(37.1) strob lights not workinohvd leak
```

```
bag =
bagOfWords with 913 words and 617 documents
```

	preventative	maintenance	service
	1	1	1
	0	0	0
...			

```
mdl =
```

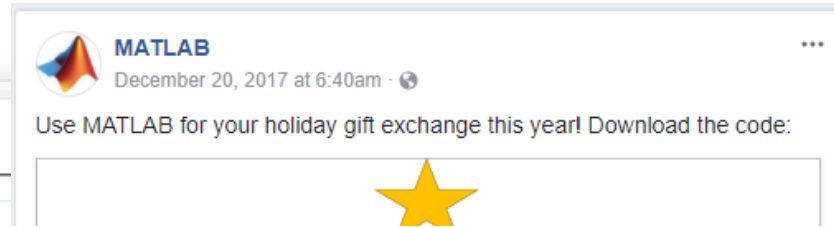
```
LdaModel with properties:
```

```
    NumTopics: 4
  TopicConcentration: 1.3262
  TopicProbabilities: [4x1 double]
  WordConcentration: 1
  Vocabulary: [1x913 string]
```

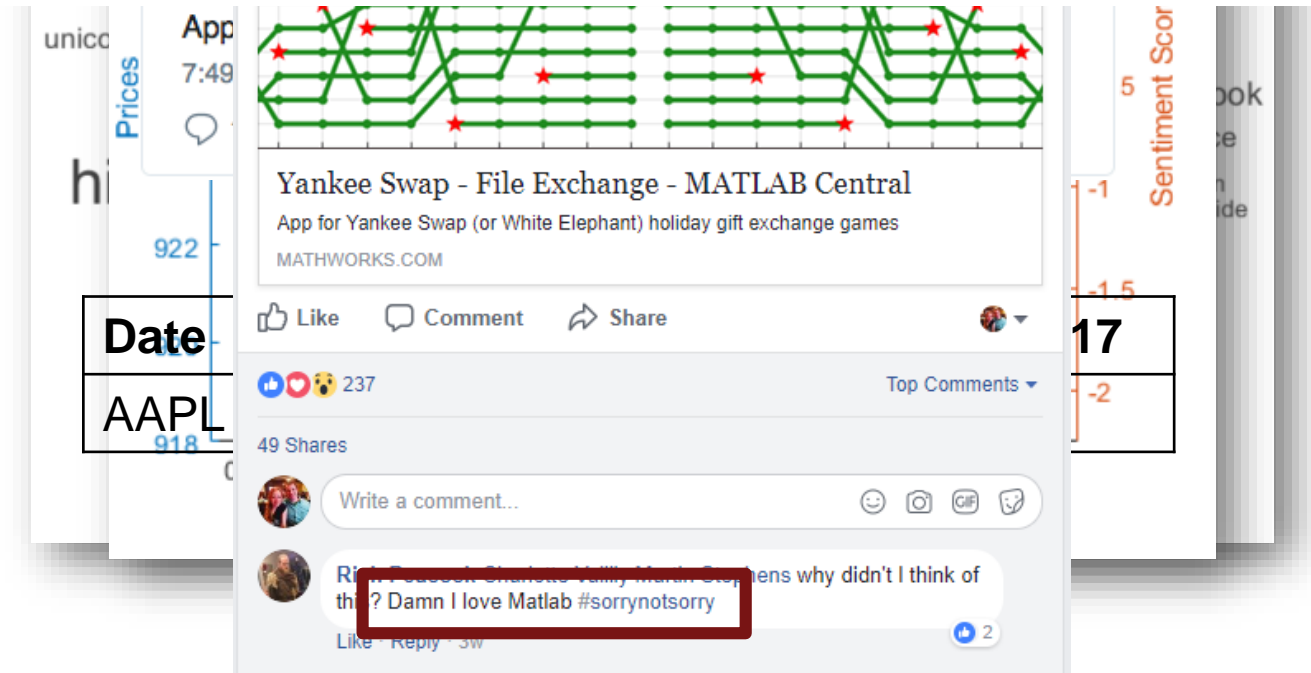
## ■ Applications

- **Sentiment Analysis: Determine if news about a product is positive/negative**
- Maintenance: Identify hidden groups of issues in maintenance logs
- Document Classification: Tag unread documents (eg. for triaging, routing, etc.)

# What is Sentiment Analysis?



Damn I love Matlab #sorrynotsorry







# Demo: Workflow

## Tweets

```
ans = 508x1 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, $O
"REPORT: Apple to build data center in Iowa: https://t.co/jwHE
"RT @theflynews: REPORT: Apple to build data center in Iowa: h
```

Preprocess text

Test data

Trained model



Score

Word Embedding

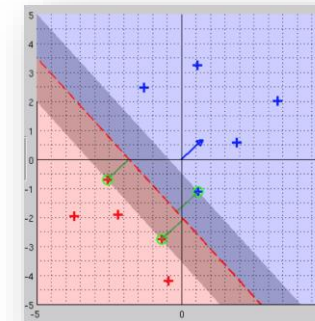
```
wordEmbedding with properties:
Dimension: 100
Vocabulary: [1x1193514 string]
```

Positive + Negative Word List

pos		neg	
str	1	str	1
1	a+	1	2-faced
2	abound	2	2-faces
3	abounds	3	abnormal
4	abundance	4	abolish
5	abundant	5	abominable
6	accessible	6	abominably
7	accessible	7	abominate

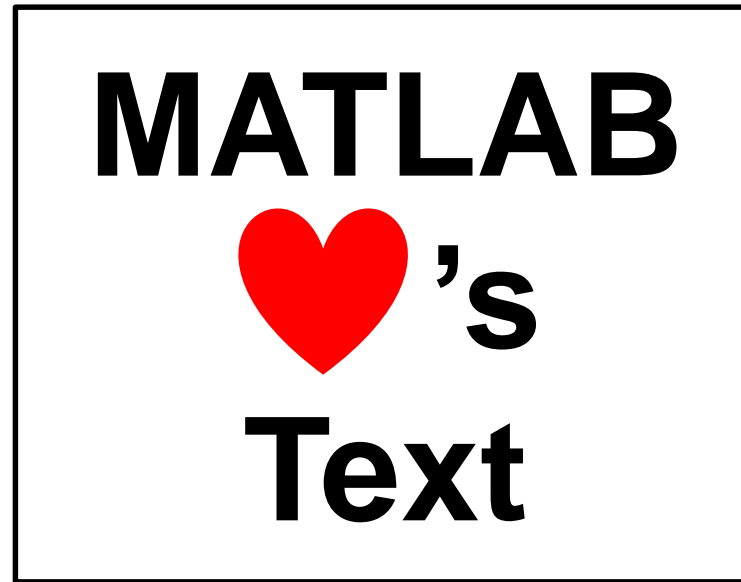
Training data

Machine Learning Model



# Sentiment Analysis Demo

## Key Takeaways



- Text data is everywhere, and contains valuable information
- Text Analytics Toolbox has tools to help you extract the signal from the noise
- Combine text with other data sources to take advantage of all your data