# From Simple Models to Topic Models

## - get a big picture of textual data

# Outline

- Create Simple Text Model for Classification
  - fitcecoc
- Analyze Text Data Using Multiword Phrases
  - bagOfNgrams
- Analyze Text Data Using Topic Models
  - fitlda

# Create Simple Text Model for Classification

- We show how to train a simple text classifier on word frequency counts using a bag-of-words model.

- We create a simple classification model which uses word frequency counts as predictors. This example trains a simple classification model to predict the category of factory reports using text descriptions.

# Load and Extract Text Data

- Load the example data. The file factoryReports.csv contains factory reports, including a text description and categorical labels for each report.

filename = "factoryReports.csv";

data = readtable(filename,'TextType','string');

head(data)

# Load and Extract Text Data

head(data)

```
ans =

8×5 table
```

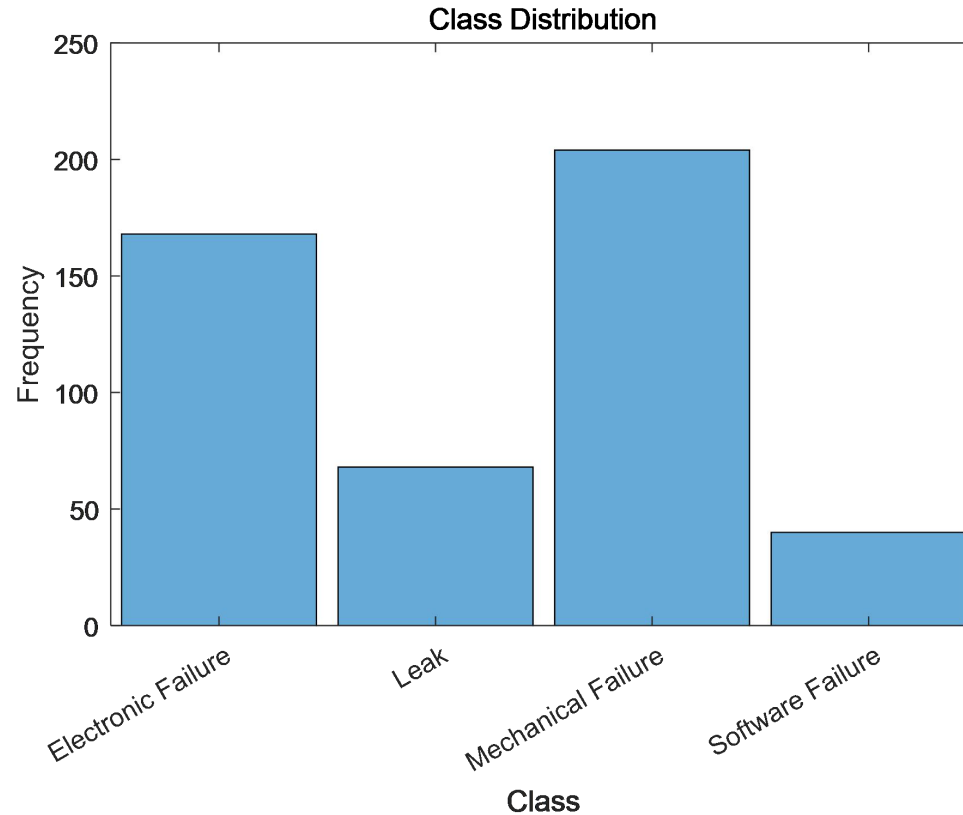| Description | Category | Urgency | Resolution | Cost |
| --- | --- | --- | --- | --- |
| "Items are occasionally getting stuck in the scanner spools." | "Mechanical Failure" | "Medium" | "Readjust Machine" | 45 |
| "Loud rattling and banging sounds are coming from assembler pistons." | "Mechanical Failure" | "Medium" | "Readjust Machine" | 35 |
| "There are cuts to the power when starting the plant." | "Electronic Failure" | "High" | "Full Replacement" | 16200 |
| "Fried capacitors in the assembler." | "Electronic Failure" | "High" | "Replace Components" | 352 |
| "Mixer tripped the fuses." | "Electronic Failure" | "Low" | "Add to Watch List" | 55 |
| "Burst pipe in the constructing agent is spraying coolant." | "Leak" | "High" | "Replace Components" | 371 |
| "A fuse is blown in the mixer." | "Electronic Failure" | "Low" | "Replace Components" | 441 |
| "Things continue to tumble off of the belt." | "Mechanical Failure" | "Low" | "Readjust Machine" | 38 |

# Load and Extract Text Data

- Convert the labels in the Category column of the table to categorical and view the distribution of the classes in the data using a histogram.

data.Category = categorical(data.Category);

figure; histogram(data.Category)

xlabel("Class")

ylabel("Frequency")

title("Class Distribution")

# Load and Extract Text Data



Class Distribution

# Load and Extract Text Data

- Partition the data into a training partition and a held-out test set. Specify the holdout percentage to be 10%.

cvp = cvpartition(data.Category,'Holdout',0.1);

dataTrain = data(cvp.training,:);

dataTest = data(cvp.test,:);

# Load and Extract Text Data

- Extract the text data and labels from the tables.

textDataTrain = dataTrain.Description;

textDataTest = dataTest.Description;

YTrain = dataTrain.Category;

YTest = dataTest.Category;

# Prepare Text Data for Analysis

- Create a function which tokenizes and preprocesses the text data so it can be used for analysis. The function preprocessText, performs the following steps in order:

1 Tokenize the text using tokenizedDocument.

2 Remove a list of stop words ("and", "of") using removeStopWords.

3 Lemmatize the words using normalizeWords.

4 Erase punctuation using erasePunctuation.

5 Remove words with 2 or fewer characters and with 15 or more characters using removeShortWords and removeLongWords.

# Prepare Text Data for Analysis

- Use the example preprocessing function preprocessText to prepare the text data.

documents = preprocessText(textDataTrain);

documents(1:5)

ans =

 5×1 tokenizedDocument:

 6 tokens: items occasionally get stuck scanner spool

 7 tokens: loud rattle bang sound come assembler piston

 4 tokens: cut power start plant

 3 tokens: fry capacitor assembler

 3 tokens: mixer trip fuse

# Prepare Text Data for Analysis

- Create a bag-of-words model from the tokenized documents.

bag = bagOfWords(documents)

bag =

 bagOfWords with properties:

 Counts: [432$\times$336 double]

 Vocabulary: [1$\times$336 string]

 NumWords: 336

 NumDocuments: 432

# Prepare Text Data for Analysis

- Remove words from the bag-of-words model that do not appear more than two times in total. Remove any documents containing no words from the bag-of-words model, and remove the corresponding entries in labels.

bag = removeInfrequentWords(bag,2);

[bag,idx] = removeEmptyDocuments(bag);

YTrain(idx) = [];

# Train Supervised Classifier

- Train a supervised classification model using the word frequency counts from the bag-of-words model and the labels.

- Train a multiclass linear classification model using fitcecoc. Specify the Counts property of the bag-of-words model to be the predictors, and the event type labels to be the response. Specify the learners to be linear. These learners support sparse data input.

XTrain = bag.Counts;

mdl = fitcecoc(XTrain,YTrain,'Learners','linear')

# Error-Correcting Output Codes

| Class | $f_0$ | $f_1$ | $f_2$ | $f_3$ | $f_4$ | $f_5$ | $f_6$ | $f_7$ | $f_8$ | $f_9$ | $f_{10}$ | $f_{11}$ | $f_{12}$ | $f_{13}$ | $f_{14}$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | 1 | 0 | 1 |
| 1 | 0 | 0 | 1 | 1 | 1 | 1 | 0 | 1 | 0 | 1 | 1 | 0 | 0 | 1 | 0 |
| 2 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 0 | 1 | 0 | 1 |
| 3 | 0 | 0 | 1 | 1 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 1 |
| 4 | 1 | 1 | 1 | 0 | 1 | 0 | 1 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 1 |
| 5 | 0 | 1 | 0 | 0 | 1 | 1 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 1 |
| 6 | 1 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 1 |
| 7 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 0 | 1 | 0 | 1 | 1 | 0 | 0 | 1 |
| 8 | 1 | 1 | 0 | 1 | 0 | 1 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 1 |
| 9 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | 1 |

Figure 1: A 15 bit error-correcting output code for a ten-class problem

# Test Classifier

- Predict the labels of the test data using the trained model and calculate the classification accuracy. The classification accuracy is the proportion of the labels that the model predicts correctly.

- Preprocess the test data using the same preprocessing steps as the training data. Encode the resulting test documents as a matrix of word frequency counts according to the bag-of-words model.

```
documentsTest = preprocessText(textDataTest);
XTest = encode(bag,documentsTest);
```

# Test Classifier

- Predict the labels of the test data using the trained model and calculate the classification accuracy.

YPred = predict(mdl,XTest);

acc = sum(YPred == YTest)/numel(YTest)

acc =

  0.8958

# Predict Using New Data

- Classify the event type of new factory reports.

str = [ "Coolant is pooling underneath sorter."

"Sorter blows fuses at start up."

"There are very loud rattling sounds coming from the assembler."];

documentsNew = preprocessText(str);

XNew = encode(bag,documentsNew);

labelsNew = predict(mdl,XNew)

labelsNew = 3$\times$1 categorical

 Leak

 Electronic Failure

 Mechanical Failure

# Analyze Text Data Using Multiword Phrases

- This example shows how to analyze text using n-gram frequency counts.

- An n-gram is a tuple of n consecutive words. For example, a bigram (the case when n = 2) is a pair of consecutive words such as "heavy rainfall". A unigram (the case when n = 1) is a single word. A bag-of-n-grams model records the number of times that different n-grams appear in document collections.

# Analyze Text Data Using Multiword Phrases

- Using a bag-of-n-grams model, you can retain more information on word ordering in the original text data. For example, a bag-of-n-grams model is better suited for capturing short phrases which appear in the text, such as "heavy rainfall" and "thunderstorm winds".

- To create a bag-of-n-grams model, use **bagOfNgrams**. You can input bagOfNgrams objects into other Text Analytics Toolbox functions such as wordcloud and fitlda.

# Load and Extract Text Data

- Load the example data. The file factoryReports.csv contains factory reports, including a text description and categorical labels for each event. Remove the rows with empty reports.

filename = "factoryReports.csv";

data = readtable(filename,'TextType','String');

# Load and Extract Text Data

- Extract the text data from the table and view the first few reports.

textData = data.Description;

textData(1:5)

ans = 5×1 string

 "Items are occasionally getting stuck in the scanner spools."

 "Loud rattling and banging sounds are coming from assembler pistons."

 "There are cuts to the power when starting the plant."

 "Fried capacitors in the assembler."

 "Mixer tripped the fuses."

# Load and Extract Text Data

- Use the example preprocessing function preprocessTest to prepare the text data.

documents = preprocessText(textData);

documents(1:5)

ans =

5×1 tokenizedDocument:

6 tokens: item occasionally get stuck scanner spool

7 tokens: loud rattle bang sound come assembler piston

4 tokens: cut power start plant

3 tokens: fry capacitor assembler ...

# Create Word Cloud of Bigrams

- Create a word cloud of bigrams by first creating a bag-of-n-grams model using bagOfNgrams.

bag = bagOfNgrams(documents)

bag =

bagOfNgrams with properties:

Counts: [480×941 double]

Vocabulary: [1×351 string]

Ngrams: [941×2 string]

NgramLengths: 2

NumNgrams: 941    NumDocuments: 480

# Create Word Cloud of Bigrams

- Visualize the bag-of-n-grams model using a word cloud.

figure

wordcloud(bag);

# Fit Topic Model to Bag-of-N-Grams

- A Latent Dirichlet Allocation (LDA) model is a topic model which discovers underlying topics in a collection of documents and infers the word probabilities in topics.

- Create an LDA topic model with 10 topics using fitlda. The function fits an LDA model by treating the n-grams as single words.

mdl = fitlda(bag,10,'Verbose',0);

# More About Latent Dirichlet Allocation

- A latent Dirichlet allocation (LDA) model is a document topic model which discovers underlying topics in a collection of documents and infers word probabilities in topics. LDA models a collection of D documents as topic mixtures $\theta_1$, …, $\theta_D$, over K topics characterized by vectors of word probabilities $\phi_1$, …, $\phi_K$. It assumes that the topic mixtures $\theta_1$, …, $\theta_D$, and the topics $\phi_1$, …, $\phi_K$ follow a Dirichlet distribution with concentration parameters $\alpha$ and $\beta$ respectively.

# More About Latent Dirichlet Allocation

- The topic mixtures $\theta_1, \ldots, \theta_D$ are probability vectors of length K, where K is the number of topics. The entry $\theta_{di}$ is the probability of topic i appearing in the dth document.

- The topics $\phi_1, \ldots, \phi_K$ are probability vectors of length V, where V is the number of words in the vocabulary. The entry $\phi_{iv}$ corresponds to the probability of the vth word of the vocabulary appearing in the ith topic.

# More About Latent Dirichlet Allocation

- Given the topics $\phi_1$, …, $\phi_K$ and Dirichlet prior α on the topic mixtures, LDA assumes the following generative process for a document:

1 Sample a topic mixture θ ～ Dirichlet(α). The random variable θ is a probability vector of length K, where K is the number of topics.

# More About Latent Dirichlet Allocation

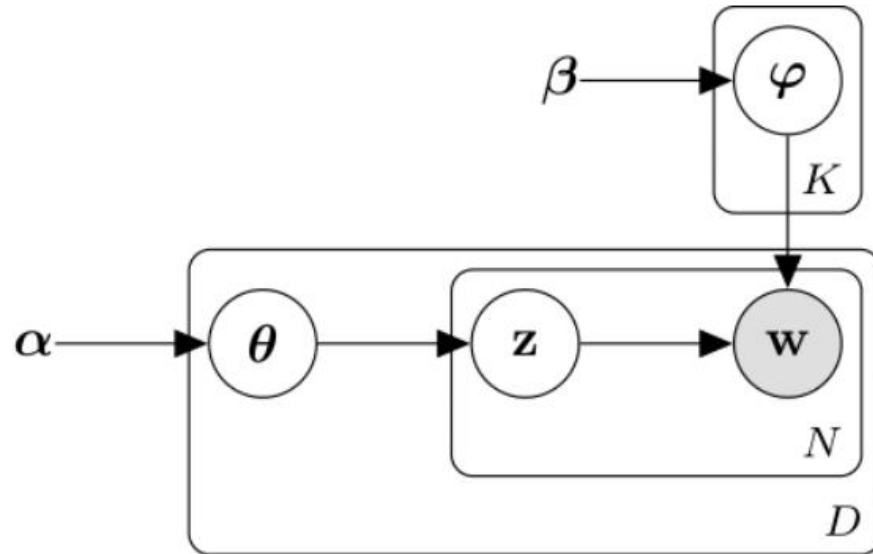2 For each word in the document:

  a) Sample a topic index $z \sim$ Categorical($\theta$). The random variable $z$ is an integer from 1 through K, where K is the number of topics.

  b) Sample a word $w \sim$ Categorical($\phi_z$). The random variable $w$ is an integer from 1 through V, where V is the number of words in the vocabulary, and represents the corresponding word in the vocabulary.

# More About Latent Dirichlet Allocation

LDA model is a probabilistic graphical model. The shaded nodes are observed variables, unshaded nodes are latent variables, nodes without outlines are model parameters.

# Fit Topic Model to Bag-of-N-Grams

- Visualize the first four topics as word clouds. The word clouds highlight commonly co-occurring bigrams in the LDA topics. The function plots the bigrams with sizes according to their probabilities for the specified LDA topics.

```
figure
for i = 1:4
  subplot(2,2,i)
  wordcloud(mdl,i);
  title("LDA Topic " + i)
end
```

LDA Topic 1

LDA Topic 2

LDA Topic 3

LDA Topic 4

# Analyze Text Using Longer Phrases

- To analyze text using longer phrases, specify the 'NGramLengths' option in bagOfNgrams to be a larger value.

- When working with longer phrases, it can be useful to keep stop words in the model. For example, to detect the phrase "is not happy", keep the stop words "is" and "not" in the model.

- Preprocess the text. Erase the punctuation using erasePunctuation, and tokenize using tokenizedDocument.

cleanTextData = erasePunctuation(textData);

documents = tokenizedDocument(cleanTextData);

# Analyze Text Using Longer Phrases

- To count the n-grams of length 3 (trigrams), use bagOfNgrams and specify 'NGramLengths' to be 3.

bag = bagOfNgrams(documents,'NGramLengths',3);

- Visualize the bag-of-n-grams model using a word cloud. The word cloud of trigrams better shows the context of the individual words.

figure

wordcloud(bag);

title("Text Data: Trigrams")

the constructing agent

stuck in scanner

Things stuck in

sounds heard inside

rattling sound heard

connect to the

neglects to interface

is powered on

the floor under

slight signs of

of the robot

over the place

from the mixer

is blown in

cracks appearing in

heard inside assembler

side constructing agent

appearing in the

is spraying coolant

in scanner spools

Fuse is blown

from the conveyor

of the sorting

all over the

in constructi

stuck in a

emitted from the

not move from

the robot arm

und heard in

heard inside the

starting to crack

fails to connect

the assembler

on the floor

from time to

is stuck in

in the assembler

blown in the

inside the mixer

time to time

sound in the

in the controller

software fails to

Shrill cry from

is hot to

falling off

move from time

in the bottom

constructing agent is

arm power supply

to the touch

to connect to

products are cracked

coming from assembler

in the mixer

by the robot

signs of wear

starting t

off the belt

of the mixer

pipe in the

Robot arm is

will not move

from the scanner

Some of the

to the ne

Reel of the

in the scanner

highpitched sound coming

e sorting plant

Fuse blown in

of the scanner

g sound emitted

sound heard inside

stuck in the

the construction agent

sounds from the

to interface with

heard in the

Assembler is overheating

by the assembler

caused by the

Mixer is hot

the bottom of

Burst pipe in

sounds made by

The mixer is

sound coming from

Robot arm power

hot to the

power cuts when

coolant in the

electrical sounds from

Slight ticking sound

bottom of the

from constructing agent

sound emitted from

# Analyze Text Using Longer Phrases

- View the top 10 trigrams and their frequency counts using topkngrams.

tbl = topkngrams(bag,10)

tbl=10 $\times$ 3 table

Ngram Count NgramLength

"in" "the" "mixer" 14 3

"in" "the" "scanner" 13 3

"blown" "in" "the" 9 3

"the" "robot" "arm" 7 3

"stuck" "in" "the" 6 3

"is" "spraying" "coolant" 6 3

"from" "time" "to" 6 3

"time" "to" "time" 6 3

"heard" "in" "the" 6 3

"on" "the" "floor" 6 3

# Analyze Text Data Using Topic Models

- This example shows how to use the Latent Dirichlet Allocation (LDA) topic model to analyze text data.

- A Latent Dirichlet Allocation (LDA) model is a topic model which discovers underlying topics in a collection of documents and infers the word probabilities in topics.

# Load and Extract Text Data

- Load the example data. The file factoryReports.csv contains factory reports, including a text description and categorical labels for each event.

data = readtable("factoryReports.csv",'TextType','string');

head(data)

ans=8 × 5 table

 Description Category Urgency Resolution Cost

 "Items are occasionally getting stuck in the scanner spools." "Mechanical Failure" "Medium" "Readjust Machine" 45

 "Loud rattling and banging sounds are coming from assembler pistons." "Mechanical Failure" "Medium" "Readjust Machine" 35 ...

# Load and Extract Text Data

- Extract the text data from the field Description.

textData = data.Description;

textData(1:5)

ans = 5 $\times$ 1 string

 "Items are occasionally getting stuck in the scanner spools."

 "Loud rattling and banging sounds are coming from assembler pistons."

 "There are cuts to the power when starting the plant."

 "Fried capacitors in the assembler."

 "Mixer tripped the fuses."

# Prepare Text Data for Analysis

- Use the preprocessing function preprocessText to prepare the text data.

documents = preprocessText(textData);

documents(1:5)

ans =

5×1 tokenizedDocument:

6 tokens: items occasionally get stuck scanner spool

7 tokens: loud rattle bang sound come assembler piston

4 tokens: cut power start plant

3 tokens: fry capacitor assembler

3 tokens: mixer trip fuse

# Prepare Text Data for Analysis

- Create a bag-of-words model from the tokenized documents.

bag = bagOfWords(documents)

bag =

bagOfWords with properties:

Counts: [480$\times$351 double]

Vocabulary: [1$\times$351 string]

NumWords: 351

NumDocuments: 480

# Prepare Text Data for Analysis

- Remove words from the bag-of-words model that have do not appear more than two times in total. Remove any documents containing no words from the bag-of-words model.

bag = removeInfrequentWords(bag,2);

bag = removeEmptyDocuments(bag)

bag =

bagOfWords with properties:

Counts: [480×162 double]

Vocabulary: [1×162 string]

NumWords: 162

NumDocuments: 480

# Fit LDA Model

- Fit an LDA model with 7 topics.

numTopics = 7;

mdl = fitlda(bag,numTopics,'Verbose',0);

- If you have a large dataset, then the stochastic approximate variational Bayes solver is usually better suited as it can fit a good model in fewer passes of the data. The default solver for fitlda (collapsed Gibbs sampling) can be more accurate at the cost of taking longer to run. To use stochastic approximate variational Bayes, set the 'Solver' option to 'savb'.

# Visualize Topics Using Word Clouds

- You can use word clouds to view the words with the highest probabilities in each topic. Visualize the first four topics using word clouds.

figure;

for topicIdx = 1:4

  subplot(2,2,topicIdx)

  wordcloud(mdl,topicIdx);

  title("Topic " + topicIdx)

end

Topic 1

Topic 2

Topic 3

Topic 4

# View Mixtures of Topics in Documents

- Use transform to transform the documents into vectors of topic probabilities.

newDocument = tokenizedDocument("Coolant is pooling underneath sorter.");
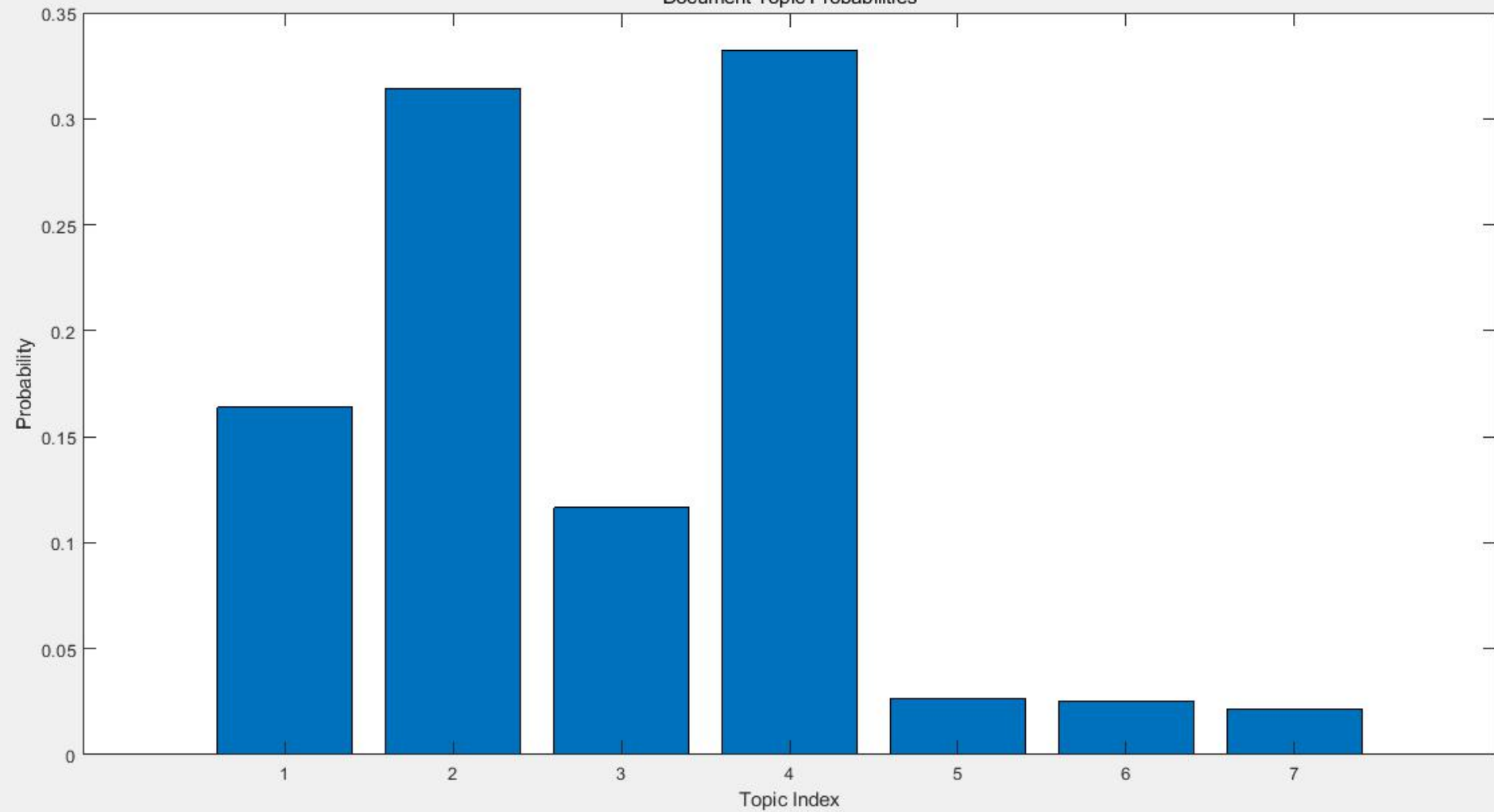
topicMixture = transform(mdl,newDocument);

figure

bar(topicMixture)

xlabel("Topic Index")

ylabel("Probability")

title("Document Topic Probabilities")

Document Topic Probabilities

# View Mixtures of Topics in Documents

- Visualize multiple topic mixtures using stacked bar charts. Visualize the topic mixtures of the first 5 input documents.

```
figure
topicMixtures = transform(mdl,documents(1:5));
barh(topicMixtures(1:5,:),'stacked')
xlim([0 1])
title("Topic Mixtures")
xlabel("Topic Probability")
ylabel("Document")
legend("Topic " + string(1:numTopics),'Location','northeastoutside')
```