

# Data Science for Business

## *Text Representation and Mining*

Asst. Prof. Teerapong Leelanupab (Ph.D.)  
Faculty of Information Technology  
King Mongkut's Institute of Technology Ladkrabang (KMITL)



Week 12

# Overview

---

- **Working with Text Data**
  - Unstructured Text
  - Tokenising Text
  - Bag-of-Words Representations
  - Measuring Similarity
  - Challenges in Text Mining
  - Pre-Processing Text with Python
- Term Weighting
- Text Classification

# Working with Text Data

# Dealing with Text

---

- Data are represented in ways natural to problems from which they were derived
- Vast amount of text..
- If we want to apply the many data mining tools that we have at our disposal, we must
  - either engineer the data representation to match the tools (**representation engineering**), or
  - build new tools to match the data

# Why Text is Difficult

---

- Text is “unstructured”
  - Linguistic structure is intended for human communication and not computers
- Word order matters sometimes
- Text can be dirty
  - People write ungrammatically, misspell words, abbreviate unpredictably, and punctuate randomly
  - Synonyms, homograms, abbreviations, etc.
- Context matters

# Unstructured Text

Most textual data arrives in an unstructured form without any pre-defined organisation or format, beyond natural language. The vocabulary, formatting, and quality of the text can vary significantly.

## United Airlines shares plummet after passenger dragged from plane

Shares plummeted Tuesday, wiping close to \$1bn off the holding company's value, after a man was violently removed from a flight by aviation police

Shares in United Airlines' parent company plummeted on Tuesday, wiping close to \$1bn off of the company's value, a day after a viral video showing police forcibly dragging a passenger off one of its plane [became a global news sensation](#).

The value of the carrier's holding company, United Continental Holdings, had fallen over 4% before noon, close to \$1bn less than the \$22.5bn as of Monday's close, according to FactSet data.

Ai woz lyin on teh stairz n [@vorvolak](#) steppd on meh! She sez she noes seez meh buh ai woz dere!  

Translated from Indonesian by  bing

Could not translate Tweet

Wrong translation?

## AUSTEN'S NOVELS EMMA

LONDON: PBINTED BY SrOTTISWOODE AND CO., NEW-STItEEI SQUASH AND PaFILIAJONT STRELT

## CHAPTER I.

MMA AVOODHOUSE, handsome, clever, and rich, with a comfortable home and happy dis position, seemed to unite some of the best blessings of existence; and had lived nearly twenty-one years in the world with very little to distress or vex her. She was the youngest of the two daughters of a most affec. tionate, indulgent father; and had, in consequence of her sister's marriage, been mistress of his house from a very early period. Her mother had died too long ago for her to have more than an indistinct remembrance of her caresses, and her place had been supplied by an excellent woman as gover ness, who had fallen little short of a mother in affection. Sixteen years had Miss Taylor been in Mr. AWoodhouse's family, less as a governess than a friend, very fond of both daughters.

Great seeing you!

Lol, wuz gr8 2 c u 2  
but omg gtg ttly!

# Common Tasks in Text Mining

---

- *Document classification*: e.g. Is a new email spam or non-spam?
- *Topic modelling*: e.g. What are the main subjects being discussed around EU Brexit negotiations?
- *Sentiment analysis*: e.g. Are Twitter users talking positively or negatively about the new iPhone?
- *Review mining*: e.g. Can we extract the most positive and negative features in reviews for a given hotel on TripAdvisor?
- *Authorship attribution*: e.g. Did Shakespeare write all his own dramas?
- *Genre classification*: e.g. Can we automatically assign a new novel to an Amazon genre category?
- *Moderation*: e.g. Can we identify potentially abusive or offensive posts on forums or on social media platforms?
- *Plagiarism detection*: e.g. Are two submitted reports unusually similar?

# Tokenising Text

ตัดคำ

- Raw text documents are textual, not numeric. The first step in analysing unstructured documents is to split the raw text into individual **tokens**, each corresponding to a single **term** (word).
- In the simplest case we might split a text string by whitespace:

```
text = "IMF cuts global growth outlook"  
text.split(" ")
```

1 เทอม = 1 token

```
[ 'IMF', 'cuts', 'global', 'growth', 'outlook' ]
```

- In most cases we need to deal with other types of punctuation:

Bad-news ahead? The IMF cuts outlook, says:'growth in doubt...'

- For some types of text, certain characters can have significance:

Discover your #Career pathway with  
@gradireland! See [bit.ly/23aPLZt](https://bit.ly/23aPLZt) for pathway  
graphics & videos #reallife

# Tokenising Text

---

- Scikit-learn provides intelligent tokenisation functions for dealing with English text - i.e. converting a string to a list of tokens.

```
from sklearn.feature_extraction.text import CountVectorizer  
tokenize = CountVectorizer().build_tokenizer()
```

```
text = "IMF cuts global growth outlook"  
tokenize(text)  
['IMF', 'cuts', 'global', 'growth', 'outlook']
```

```
text = "Bad-news ahead? The IMF cuts outlook, says:Growth in doubt"  
tokenize(text)
```

```
['Bad', 'news', 'ahead', 'The', 'IMF', 'cuts', 'outlook', 'says', 'Growth', 'in', 'doubt']
```

- Punctuation and single letter tokens are removed. A standard tokeniser is not always suitable...

```
tweet = "Discover your #Career pathway with @gradireland! See http://bit.ly/23aPLZt"  
tokenize(tweet)
```

```
['Discover', 'your', 'Career', 'pathway', 'with', 'gradireland', 'See', 'http', 'bit',  
'ly', '23aPLZt']
```

# Text Representation

---

- **Goal:** Take a set of documents –each of which is a relatively free-form sequence of words– and turn it into our familiar feature-vector form
- A collection of documents is called a *corpus*
- A *document* is composed of individual *t*okens or *terms*
- *Each document is one instance*
  - *but we don't know in advance what the features will be*

# Counting Term Frequencies

- Once we have performed tokenisation, we might then count the frequency of occurrence of terms (tokens) in each document.

Growth forecasts cut as IMF warns of risk. The IMF cut its global growth forecast for 2016. It warned of widespread global stagnation risk and stated the global economy could be vulnerable in 2016.

ข้อมูลของเอกสารเดียว

Term	Frequency
global	3
risk	2
cut	2
IMF	2
2016	2
...	...

- We can repeat this process to compute frequencies across an entire corpus and the sum the frequencies.

Growth forecasts cut as IMF warns of risk. The IMF  
c IMF chief economist Maurice Obstfeld said in a  
v sta The Fund called on global policymakers attending  
g ef the IMF and World Bank meetings, being held in  
the Washington, to take coordinated actions to boost  
demand with structural economic reforms in 2016.

ข้อมูลของหลายเอกสาร (ทั้ง corpus)

Term	Frequency
global	5
imf	4
growth	3
2016	3
economy	3
...	...

# Bag-of-Words Model for Documents

- How can we go from tokens to numeric features?
- **Bag-of-words model:** Each document is represented by a vector in a  $m$ -dimensional coordinate space, where  $m$  is total number of unique terms across all documents (the corpus **vocabulary**).

## Document 1:

Forecasts cut as IMF issues warning

## Document 2:

IMF and WBG meet to discuss economy

## Document 3:

WBG issues 2016 growth warning

## Example:

When we tokenise our corpus of 3 documents, we have a vocabulary of 14 distinct terms (when no filtering is applied).

```
vocab = set()
for doc in corpus:
    tokens = tokenize(doc)
    for tok in tokens:
        vocab.add(tok)
print(vocab)

{'2016', 'Forecasts', 'IMF', 'WBG', 'and',
 'as', 'cut', 'discuss', 'economy', 'growth',
 'issues', 'meet', 'to', 'warning'}
```

ไม่ซ้ำกันเลย

2016
and
as
cut
discuss
economy
Forecasts
growth
IMF
issues
meet
to
warning
WBG

# “Bag of Words”

---

- Treat every document as just a collection of individual words
  - Ignore grammar, word order, sentence structure, and (usually) punctuation
  - Treat every word in a document as a potentially important keyword of the document
- What will be the feature's value in a given document?
  - Each document is represented by a one (if the token is present in the document) or a zero (the token is not present in the document)
- Straightforward representation
- Inexpensive to generate
- Tends to work well for many tasks

# Bag-of-Words Representation

- Each document can be represented as a **term vector**, with an entry indicating the number of time a term appears in the document:

Document 1:

Forecasts cut as IMF issues  
warning

2016	Forecasts	IMF	WBG	and	as	cut	discuss	economy	growth	issues	meet	to	warning
0	1	1	0	0	1	1	0	0	0	1	0	0	1

- By transforming all documents in this way, and stacking them in rows, we create a full **document-term matrix**:

ใช้ term frequency ที่ไม่ถูก normalize  
ในการเก็บ

Document 2:

IMF and WBG meet to  
discuss economy

2016	Forecasts	IMF	WBG	and	as	cut	discuss	economy	growth	issues	meet	to	warning
0	1	1	0	0	1	1	0	0	0	1	0	0	1
0	0	1	1	1	0	0	1	1	0	0	1	1	0
2	0	0	1	0	0	0	0	0	1	1	0	0	1

Document 3:

2016: WBG issues 2016  
growth warning

3 Documents x 14 Terms

# Bag-of-Words Representation

---

- The position (context) of terms within the original document text is lost when using this model.

2016	Forecasts	IMF	WBG	and	as	cut	discuss	economy	growth	issues	meet	to	warning
0	1	1	0	0	1	1	0	0	0	1	0	0	1

- The size of the vocabulary is often large, meaning that the resulting vectors can be very **high dimensional**.
- Fortunately, most values in the document-term matrix will be zeros, since for a given document less than a couple thousands of distinct terms will be used.
- We say bag-of-word models are typically high-dimensional **sparse** datasets. Document-term matrix often has > 98% zero values.

# Bag-of-Words in Python

- Scikit-learn includes functionality to easily transform a collection of strings containing documents into a document-term matrix.

Our input, `documents`, is a list of strings. Each string is a separate document.

```
from sklearn.feature_extraction.text import CountVectorizer
vectorizer = CountVectorizer()
X = vectorizer.fit_transform(documents)
```

Our output, `X`, is a sparse NumPy 2D array with rows corresponding to documents and columns corresponding to terms.

- Once the matrix has been created, we can access the list of all terms and an associated dictionary (`vocabulary_`) which maps each unique term to a corresponding column in the matrix.

```
terms = vectorizer.get_feature_names()
len(terms)
```

3288

How many terms in the vocabulary?

```
vocab = vectorizer.vocabulary_
vocab["world"]
```

3246

Which column corresponds to term?

หมายเลข index ที่พ่อให้

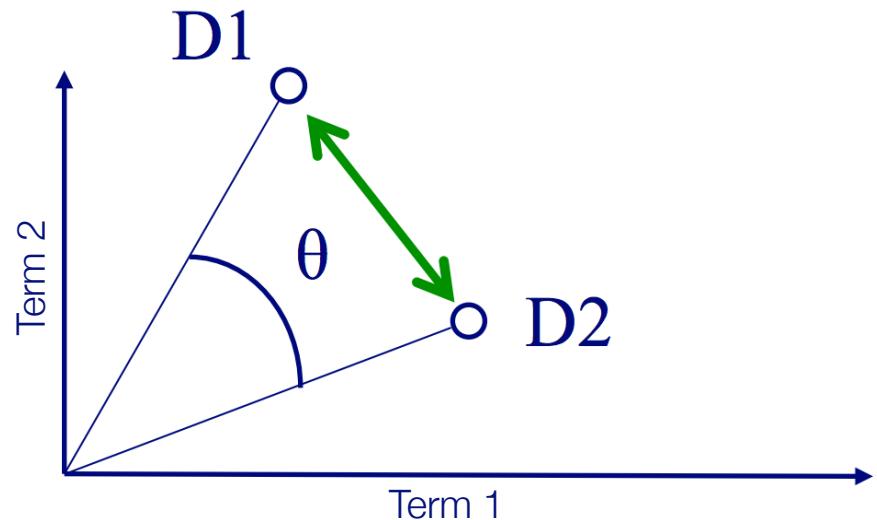
# Measuring Similarity

---

- **Cosine similarity:** Most common approach for measuring similarity between two documents in a bag-of-words representation is to look at the cosine of the angle between their term vectors.
- Motivation: vectors for documents containing similar terms will point in the same direction in the  $m$ -dimensional vector space.

$$\cos(D_1, D_2) = \frac{D_1 \cdot D_2}{\|D_1\| \|D_2\|}$$

$$\text{where } \|D\| = \sqrt{\sum_{i=1}^m d_i^2}$$



- Cosine similarity score is 1 if two documents are identical, 0 if two documents share no terms in common.

# Measuring Similarity

Document 1:

Forecasts cut as IMF issues warning

Document 2:

IMF and WBG meet to discuss economy

Document 3:

IMF and WBG meet to discuss growth

and	as	cut	discuss	economy	Forecasts	growth	IMF	issues	meet	to	warning	WBG
0	1	1	0	0	1	0	1	1	0	0	1	0
1	0	0	1	1	0	0	1	0	1	1	0	1
1	0	0	1	0	0	1	1	0	1	1	0	1

$$\cos(D1, D2) = 0.15 \quad \cos(D1, D3) = 0.15$$

$$\cos(D2, D3) = 0.86$$

We can perform the same calculations in Python...

```
from sklearn.metrics.pairwise import cosine_similarity
print( "cos(D1,D2) = %.2f" % cosine_similarity( X[0], X[1] ) )
print( "cos(D1,D3) = %.2f" % cosine_similarity( X[0], X[2] ) )
print( "cos(D2,D3) = %.2f" % cosine_similarity( X[1], X[2] ) )
```

`cos(D1,D2) = 0.15`

`cos(D1,D3) = 0.15`

`cos(D2,D3) = 0.86`

Measure similarity between rows of the document-term matrix X

eg อาจมี query ให้เป็น 1 เอกสารแล้วดูว่าเมื่อ user comment มา  
เราจะได้ไปเทียบว่าใกล้เคียงกับเอกสารไหนมากที่สุด

# Challenges in Text Mining

---

- **The Sparsity Problem:** Natural language use often involves huge vocabularies, so most documents will share very few words.
- The use of N-grams also increases sparsity, as even fewer documents will share the same phrases.
- This problem is exacerbated due to...
  - **Synonymy:** languages often use many different words to refer to identical or closely related concepts (e.g. ‘residence’, ‘abode’).
  - **Polysemy:** a single word can have multiple *related* meanings (e.g. ‘bank’ - a financial institution or the building housing it?).  
    ความหมายมากกว่า 1 แต่ไม่สัมพันธ์กันเลย
  - **Homonymy:** two or more words have the identical form or the same spelling or pronunciation, but have different unrelated meanings and origin (e.g. ‘jaguar’ - car or ‘jaguar’ - animal?) or (e.g. ‘bank’ - a financial institution or ‘bank’ - land bordering on a river?).
- This applies when analyzing a monolingual corpus. Text analysis becomes far more difficult when working with a multilingual corpus.

# Polysemy vs. Homonymy

---

Polysemy is the coexistence of many possible meanings for a word or phrase

Homonymy refers to the existence of unrelated words that look or sound the same

Has different, but related meanings

Has completely different meanings

Has related word origins

Has different origins

Polysemous words are listed under one entry in dictionaries

Homonyms are listed separately in dictionaries

Polysemous words can be understood if you know the meaning of one word

Meaning of homonyms cannot be guessed since the words have unrelated meanings

# Challenges in Text Mining

---

- The Sparsity Problem has several practical consequences...
  - Additional memory and processing requirements due to more dimensions (terms).
  - Analytical problems related to sparsity - e.g. when using a bag-of-words representation, we can fail to identify documents which are related to the same concepts.

Document 1:

jaguars are expensive cars

Document 2:

a jaguar is a costly vehicle

Document 3:

the jaguar, a feline animal

a	animal	are	cars	costly	expensive	feline	is	jaguar	jaguars	the	vehicle
0	0	1	1	0	1	0	0	0	1	0	0
1	0	0	0	1	0	0	1	1	0	0	1
1	1	0	0	0	0	1	0	1	0	1	0

$$\cos(D1, D2) = 0.0 \quad \cos(D2, D3) = 0.4$$

# Text Pre-processing

---

- The number of terms used to represent documents (and hence the sparsity) is often reduced by applying a number of simple preprocessing techniques before building a document-term matrix:
  - **Minimum term length:** Exclude terms of length < 2
  - **Case conversion (Text Normalization):** Converting all terms to lowercase.
  - **Stemming:** Process by which endings are removed from terms in order to remove things like tense or plurals:  
e.g. compute, computing, computer = comput  
ลบรูปคำให้อยู่ในรากศัพท์ของคำ
  - **Lemmatization:** Process to reduce a term to its canonical form. This is a more advanced form of stemming.
  - **Stop-word filtering/removal:** Remove terms that appear on a pre-defined filter list of terms that are highly frequent and do not convey useful information (e.g. and, the, while)  
มืออยู่ใน document เยอะเกินไป ไม่ได้ส่งไปถึงความหมายของเอกสารเลย สามารถถูกตัดออกได้
  - **Low frequency filtering:** Remove terms that appear in very few documents.

ความถี่น้อยก็ควรถูกตัดออก พอมัน rare มันจะใช้ในเอกสารอื่นน้อย  
ตามไปด้วย และถ้าไม่มีประโยชน์ มันจะไปลด similarity กับเอกสาร  
อื่นเปล่าๆ เจอก็ตัดไปเลยก็ได้

# Further Text Preprocessing

- Further preprocessing steps can be applied directly using the **CountVectorizer** class by passing appropriate parameters - e.g.:

```
from sklearn.feature_extraction.text import CountVectorizer
vectorizer = CountVectorizer(
    stop_words=custom_list,
    min_df=20,
    max_df=1000,
    lowercase=False,
    ngram_range=2)
A = vectorizer.fit_transform(documents)
```

Parameter	Explanation
stop_words=custom_list	Pass in a custom list containing terms to filter.
min_df=20	Filter those terms that appear in < 20 documents.
max_df=1000	Filter those terms that appear in > 1000 documents.
lowercase=False	Do not convert text to lowercase. Default is True.
ngram_range=2	Include phrases of length 2, instead of just single words.

# Text Preprocessing in Python

---

- By default Scikit-learn converts tokens to lowercase and removes tokens of length 1 (i.e. single letters).
- Scikit-learn allows us to perform other simple preprocessing steps by adapting the CountVectorizer.

```
vectorizer = CountVectorizer(stop_words="english")
x = vectorizer.fit_transform(documents)
```

Use built-in stop-words by specifying the language

```
mywords = [ "and", "the", "am", "pm" ]
vectorizer = CountVectorizer(stop_words=mywords)
x = vectorizer.fit_transform(documents)
```

Use a custom list of stopwords

Filter terms appearing in less than 5 documents:

```
vectorizer = CountVectorizer(min_df = 5)
x = vectorizer.fit_transform(documents)
```

Multiple preprocessing steps...

```
vectorizer = CountVectorizer(stop_words="english", min_df = 5)
x = vectorizer.fit_transform(documents)
```

# Text Preprocessing in Python

- To stem or lemmatise tokens, we need to use functionality from another text analysis library: NLTK (<http://www.nltk.org>)
- Apply stemming to individual words:

Install via command line:

```
conda install nltk
```

```
from nltk.stem.porter import PorterStemmer
words = ['cycles', 'insists', 'computer', 'flying', 'cars']
stemmer = PorterStemmer()
for w in words:
    print( stemmer.stem(w) )
cycl insist comput fli car
```

- Apply lemmatisation to individual words:

```
from nltk.stem import WordNetLemmatizer
lemmatizer = WordNetLemmatizer()
for w in words:
    print( lemmatizer.lemmatize(w) )
cycle insists computer flying car
```

# Overview

---

- **Working with Text Data**
  - Unstructured Text
  - Tokenising Text
  - Bag-of-Words Representations
  - Measuring Similarity
  - Challenges in Text Mining
  - Pre-Processing Text with Python
- **Term Weighting**
- **Text Classification**

# Term Weighting

# Notations

---

ค่าน้ำหนักของ term k ใน doc i

- $w_{i,k}$  = weight of  $k^{th}$  term in a document  $D_i$
- $f_{i,k}$  = number of occurrences of  $k^{th}$  term in a document  $D_i$   
(term frequency)
- $N$  = number of documents in the collection
- $n_k$  = number of documents containing a term  $k$

จำนวนของเอกสารที่มี term k ปรากฏเท่านั้น

# Term Weighting

---

- Two demands on the weight
  - The degree to which a particular document about a topic (or a particular term)
    - Repetition is an indication of emphasis
    - Let us call it ***aboutness<sub>k</sub>***
  - The degree to which a keyword discriminates documents in the collection
    - Let us call it ***discrim<sub>k</sub>***

เอกสารนั้น พูดเกี่ยวกับอะไร

$$w_{i,k} \propto \textbf{\textit{aboutness}}_k \times \textbf{\textit{discrim}}_k$$

คำนั้น มีส่วนที่ช่วยในการแยกแยะ  
เอกสารเท่าได

$$= tf_{i,k} \times idf_k$$

# Term Frequency

---

- Use the word count (frequency) in the document instead of just a zero or one
  - Differentiates between how many times a word is used

# Normalized Term Frequency

---

ต้องทำ เพราะมันหากไม่ทำ มันจะไม่แฟร์กับเอกสารที่มีความยาวน้อย แต่มีคำนั้นเยอะ

- Documents of various lengths
- Words of different frequencies
  - Words should not be *too common* or *too rare*
  - Both upper and lower limit on the number (or fraction) of documents in which a word may occur
  - Feature selection is often employed
- The raw term frequencies are normalized in some way,
  - such as by dividing each by the total number of words in the document
  - or the frequency of the specific term in the corpus

# Aboutness component

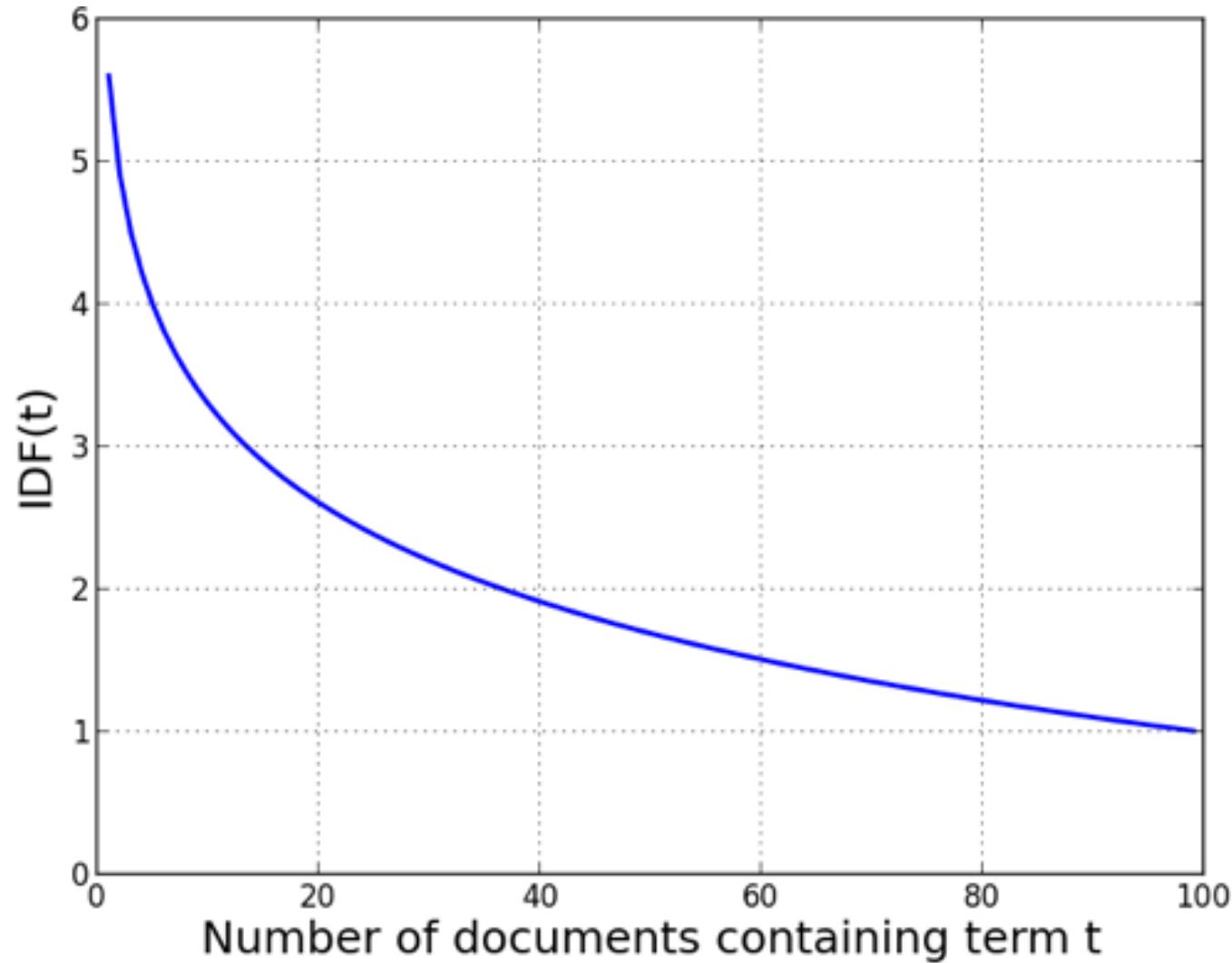
---

- **Term Frequency (tf)** reflects the importance of a term  $k$  in a document  $D_i$ ;
- It is usually normalized by the count of **all** term occurrences in a document
  - This is to reduce the effect of a long document

$$tf_{i,k} = \frac{f_{i,k}}{\sum_{j=1}^t f_{i,j}}$$

# Inverse Document Frequency (IDF)

---



# Discriminating component

---

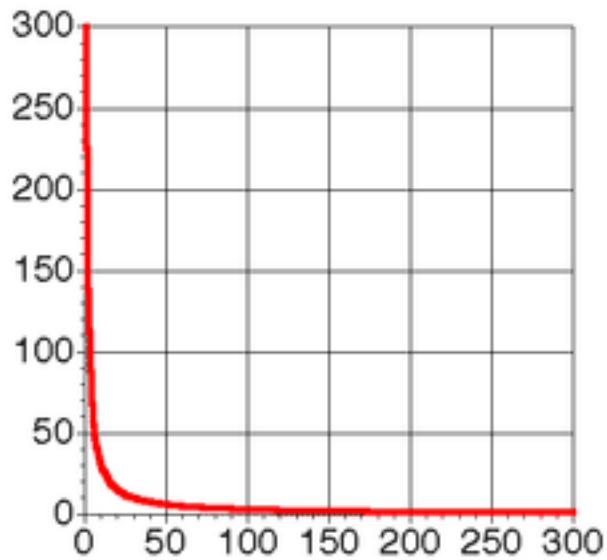
- **Inverse Document Frequency (idf)** forms a discriminating point of view
- The more documents that a term  $k$  occurs in, the less *discriminating* the term is between documents, and
- consequently, the less useful it will be in retrieval
  - It is the more specific terms that are particularly important in identifying relevant material

$$idf_k = \log \frac{N}{n_k}$$

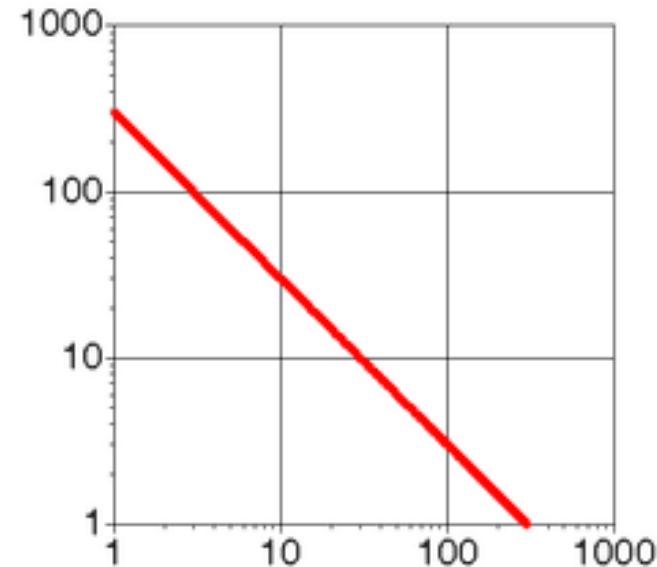
หากมีคำบางคำเกิดขึ้นในทุกเอกสาร  
แสดงว่าคำนั้น หากนำมาเป็น  
classifier จะเป็น classifier ที่ไม่ดี

# Linear and log scale

---



$$\frac{N}{D_k}$$



$$\log \frac{N}{D_k}$$

# Variants of TF-IDF weighting

$$w_{i,k} = tf_{i,k} \cdot \log \frac{N}{n_k}$$

$$w_{i,k} = tf_{i,k} \cdot \left( \log \frac{N}{n_k} + 1 \right)$$

$$w_{i,k} = tf_{i,k} \cdot \log \frac{N+1}{n_k+1}$$

$$w_{i,k} = tf_{i,k} \cdot \log \frac{(N - n_k) + 0.5}{n_k + 0.5}$$

# Term Weighting

---

- As well as including or excluding terms, we can also modify or weight the frequency values themselves.
- We can improve the usefulness of the document-term matrix by giving higher weights to more "important" terms.
- **TF-IDF**: Common approach for weighting the score for a term in a document. Several different formulations, but always consist of:
  - **Term Frequency (TF)**: Number of times a given term appears in a single document.
  - **Inverse Document Frequency (IDF)**: Function of total number of distinct documents containing a term. Effect is to penalise common terms that appear in almost every document.

Common version is  
log-based TF-IDF

$$w_{i,k} = \frac{tf_{i,k} \times (\log(\frac{N}{n_k}) + 1)}{\text{TF}} \quad \text{IDF}$$

$N$  = total number  
of documents

# Term Weighting in Scikit-learn

---

- A similar vectorisation approach can be used in Scikit-learn to produce a TF-IDF normalised document-term matrix:

```
from sklearn.feature_extraction.text import TfidfVectorizer
vectorizer = TfidfVectorizer()
A = vectorizer.fit_transform(documents)
```

The output, **A**, is a sparse NumPy array where the entries are all TF-IDF normalised.

- Again we can perform additional preprocessing steps by passing the appropriate parameter values to **TfidfVectorizer**:

```
from sklearn.feature_extraction.text import TfidfVectorizer
vectorizer = TfidfVectorizer(
    stop_words=custom_list,
    min_df=20,
    max_df=1000,
    lowercase=False,
    ngram_range=2)
A = vectorizer.fit_transform(documents)
```

# Term Weighting

---

- **Example:** In a dataset of  $N = 1000$  documents, for a document  $i$ ...

Term 'cat' appears in the document  $i$  3 times, and appears in  $n_k = 50$  documents in total.

Term 'dog' appears in the document  $i$  4 times, and appears in  $n_k = 250$  documents in total.

$$w(i, \text{cat}) = 3 \times (\log_e(1000/50) + 1) = 11.987$$

Term 'cat' is more  
"important"

$$w(i, \text{dog}) = 4 \times (\log_e(1000/250) + 1) = 9.545$$

- In Scikit-learn, we can generate a TF-IDF weighted document-term matrix by using `TfidfVectorizer` in place of `CountVectorizer`:

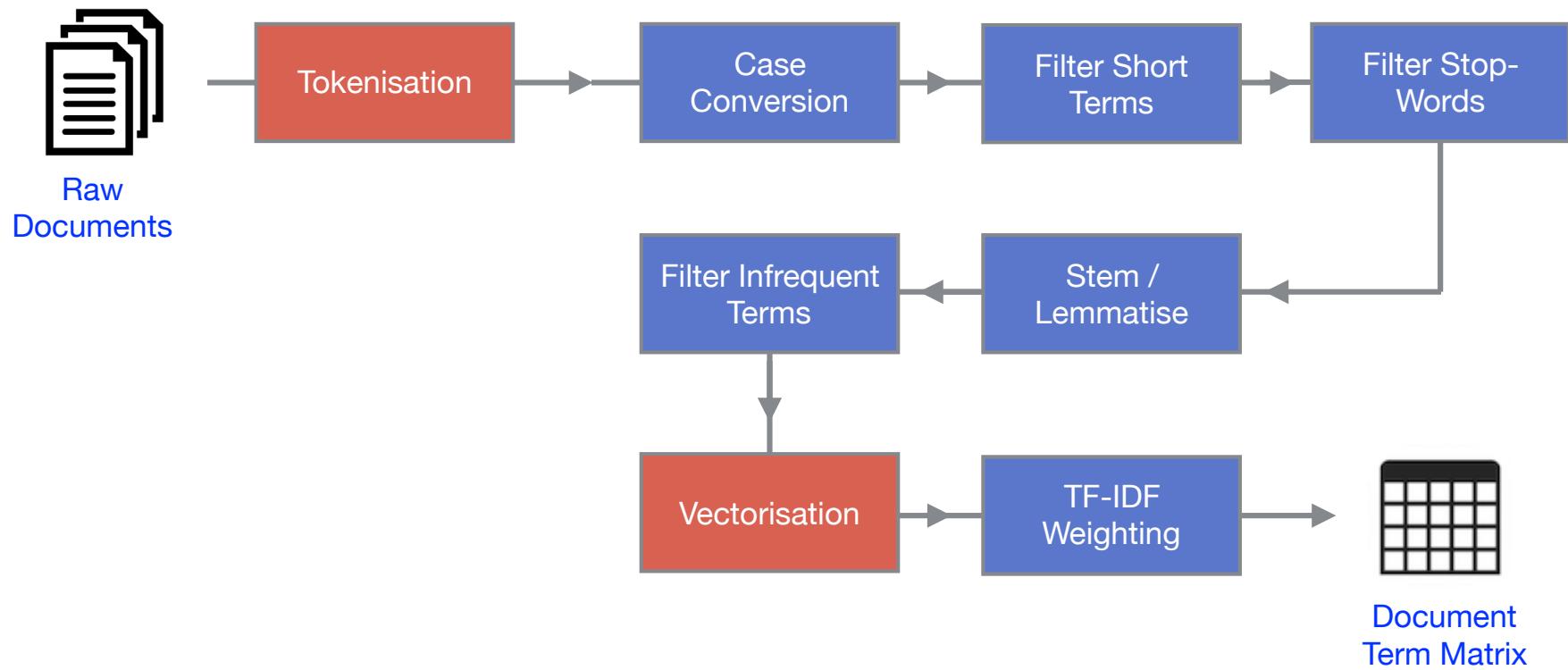
```
from sklearn.feature_extraction.text import TfidfVectorizer
vectorizer = TfidfVectorizer()
X = vectorizer.fit_transform(documents)
```

```
vectorizer = TfidfVectorizer(stop_words="english", min_df = 5)
X = vectorizer.fit_transform(documents)
```

Can include  
preprocessing steps  
as before

# Text Processing Pipeline Overview

Typical text preprocessing steps for processing a corpus...



Subsequent modelling steps can then be applied to the document-term matrix - e.g. document classification, documenter clustering.

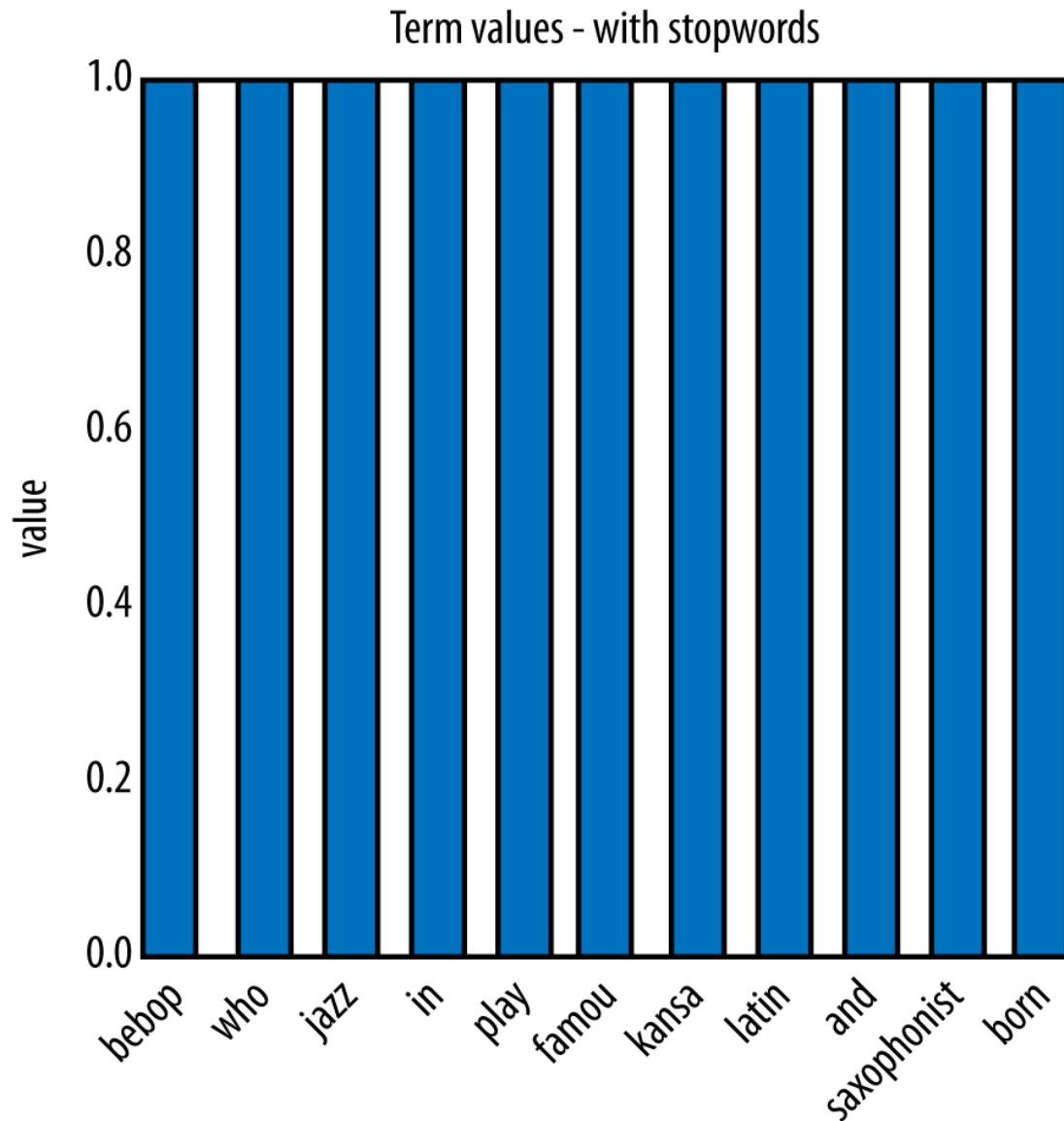
# Example: Jazz Musicians

---

- 16 prominent jazz musicians and excerpts of their biographies from Wikipedia
- Nearly 2,000 features after stemming and stop-word removal!
- Consider the sample phrase “Famous jazz saxophonist born in Kansas who played bebop and latin”

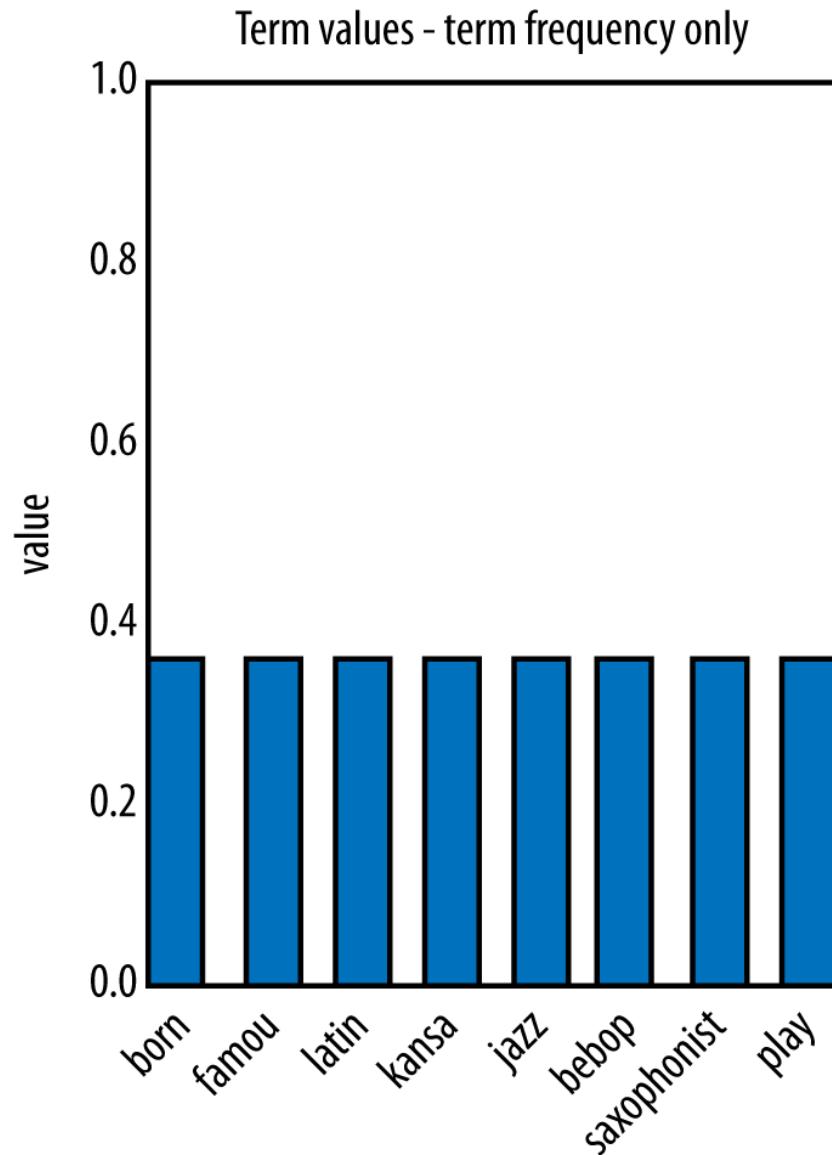
# Example: Jazz Musicians

---



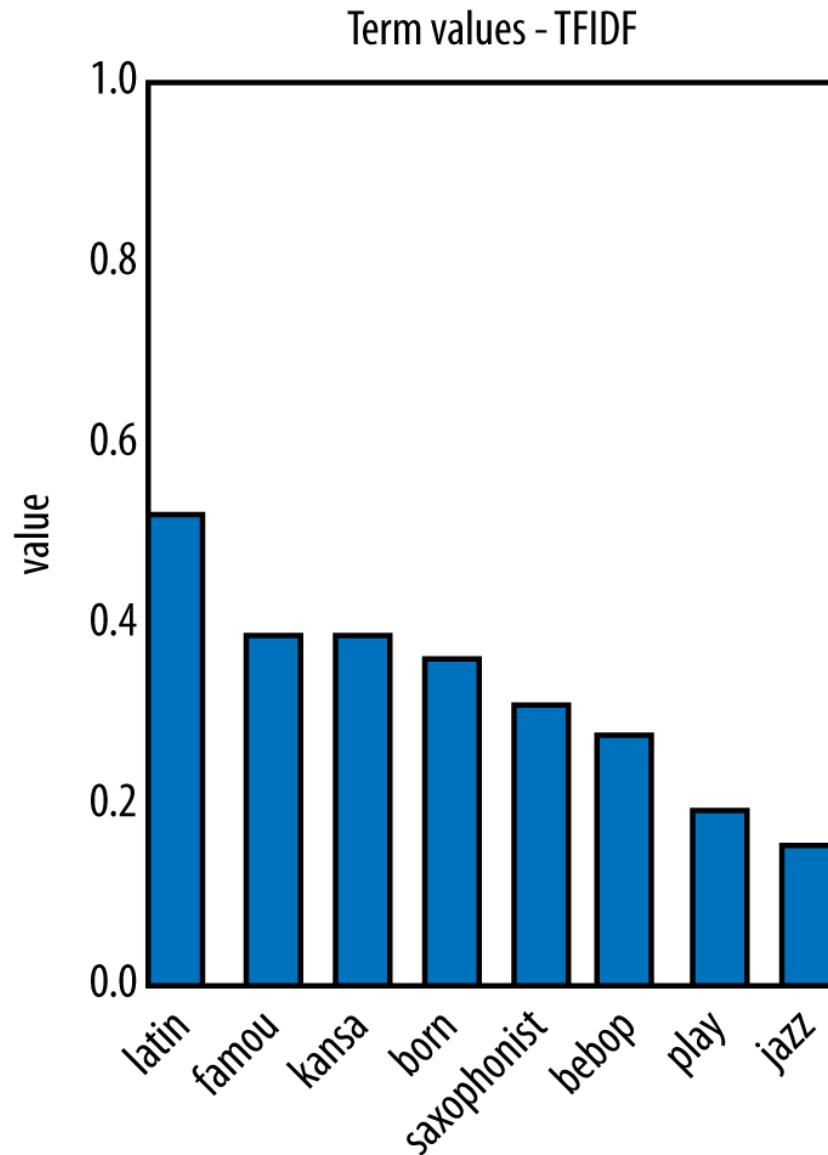
# Example: Jazz Musicians

---



# Example: Jazz Musicians

---



# Example: Jazz Musicians

Musician	Similarity	Musician	Similarity
Charlie Parker	0.135	Count Basie	0.119
Dizzie Gillespie	0.086	John Coltrane	0.079
Art Tatum	0.050	Miles Davis	0.050
Clark Terry	0.047	Sun Ra	0.030
Dave Brubeck	0.027	Nina Simone	0.026
Thelonius Monk	0.025	Fats Waller	0.020
Charles Mingus	0.019	Duke Ellington	0.017
Benny Goodman	0.016	Louis Armstrong	0.012

# Overview

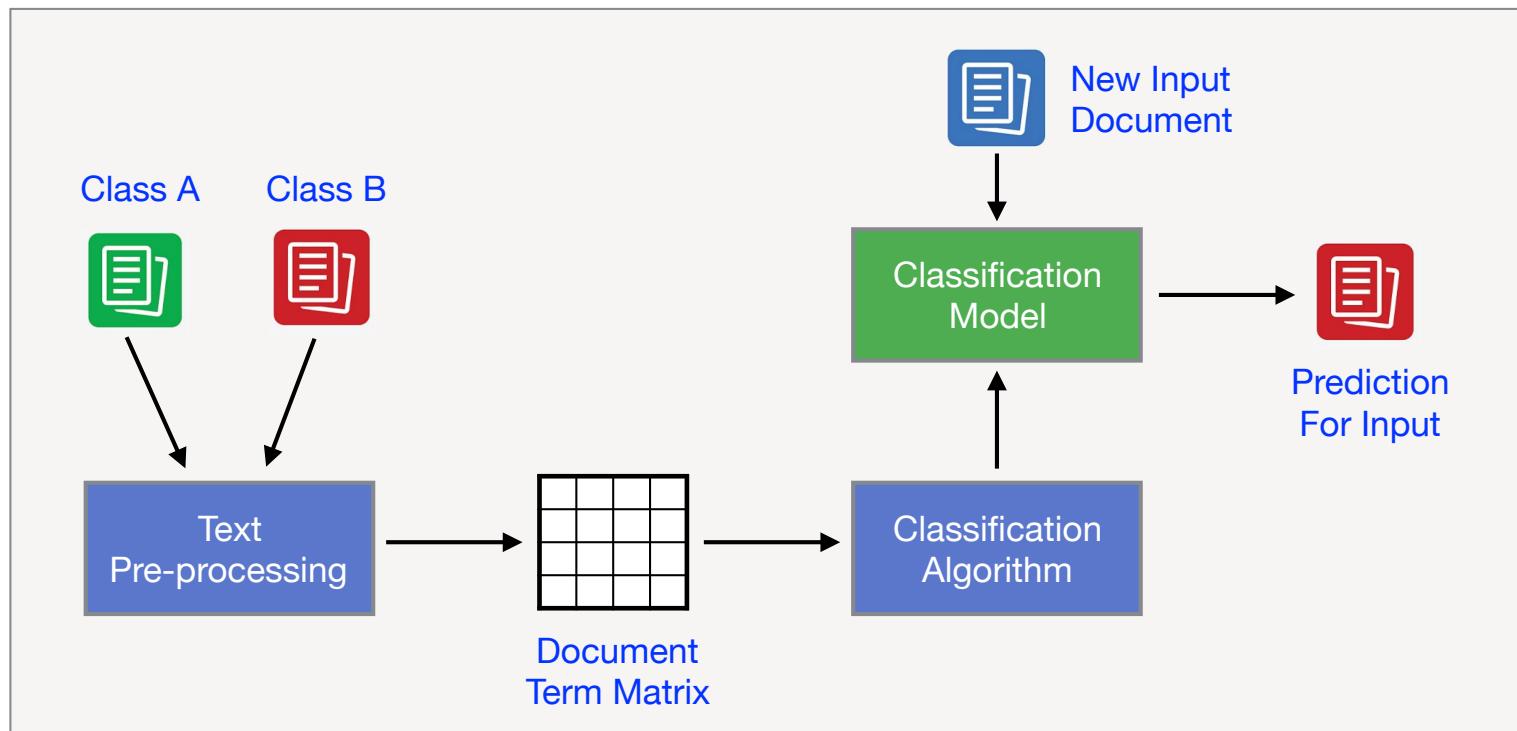
---

- **Working with Text Data**
  - Unstructured Text
  - Tokenising Text
  - Bag-of-Words Representations
  - Measuring Similarity
  - Challenges in Text Mining
  - Pre-Processing Text with Python
- **Term Weighting**
- **Text Classification**

# Text Classification

# Text Classification

- **Goal:** To learn a model from the training set so that we can accurately predict classes for new unlabeled documents.
- **Input:** Training set of labelled text documents, annotated with two or more class labels (categories).



# Text Classification

---

- When we have labelled documents to use as training data, we can apply many standard classifiers to a document-term matrix, using the functionality from Scikit-learn that we saw previously.
- As before, we will need to split training and test data for evaluations, with 2 sets of documents: `train_documents`, `test_documents`

```
vectorizer = TfidfVectorizer()  
train_X = vectorizer.fit_transform(train_documents)
```

Create document-term matrix from training documents by calling `fit_transform()`

```
from sklearn.neighbors import KNeighborsClassifier  
model = KNeighborsClassifier(n_neighbors=3)  
model.fit(train_X, train_target)
```

Build a KNN model on it

```
test_X = vectorizer.transform(test_documents)
```

Create document-term matrix from test documents. We call `transform()` to use the same vocabulary as before.

```
predicted = model.predict(test_X)
```

We can now make predictions for our test documents using the KNN model and the test document-term matrix

# Text Classification

---

- A number of general purpose classification algorithms are frequently used for classifying text documents:
  - **kNN:** Standard nearest neighbour classifier, using an appropriate similarity measure (e.g. Cosine).
  - **Naive Bayes:** Classification based on term frequency counts. Incorrectly assumes all terms are independent, but can still be effective in practice.
  - **Support Vector Machines:** Often apply SVMs with a linear kernel to calculate document similarity.
- To compare the performance of different algorithms and/or different parameter settings on the same corpus, we use standard classifier evaluation methods - e.g. measure each classifier's mean accuracy in a k-fold cross-validation experiment.

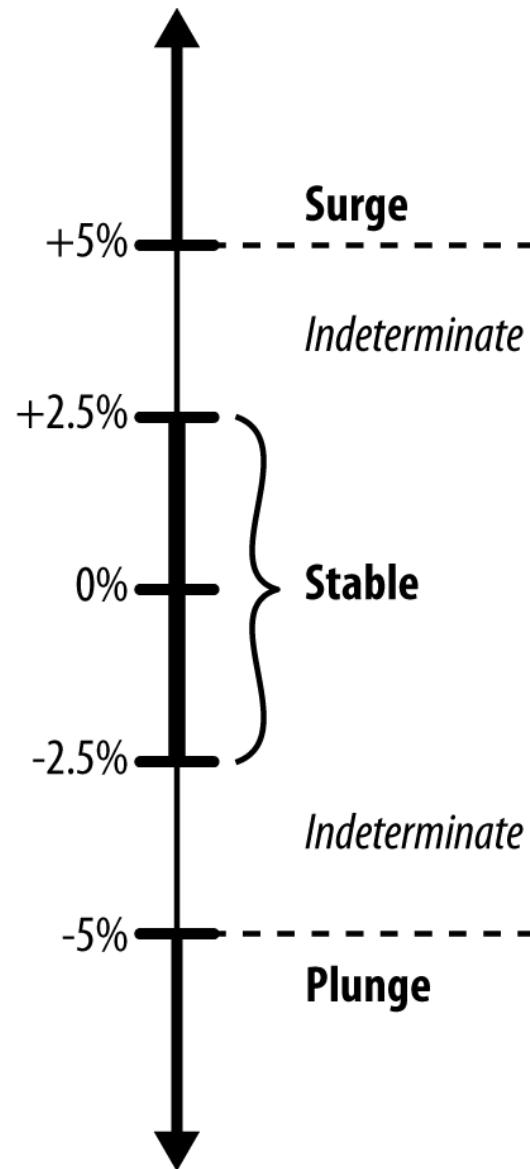
# Text Classification – Movement (Example)

---

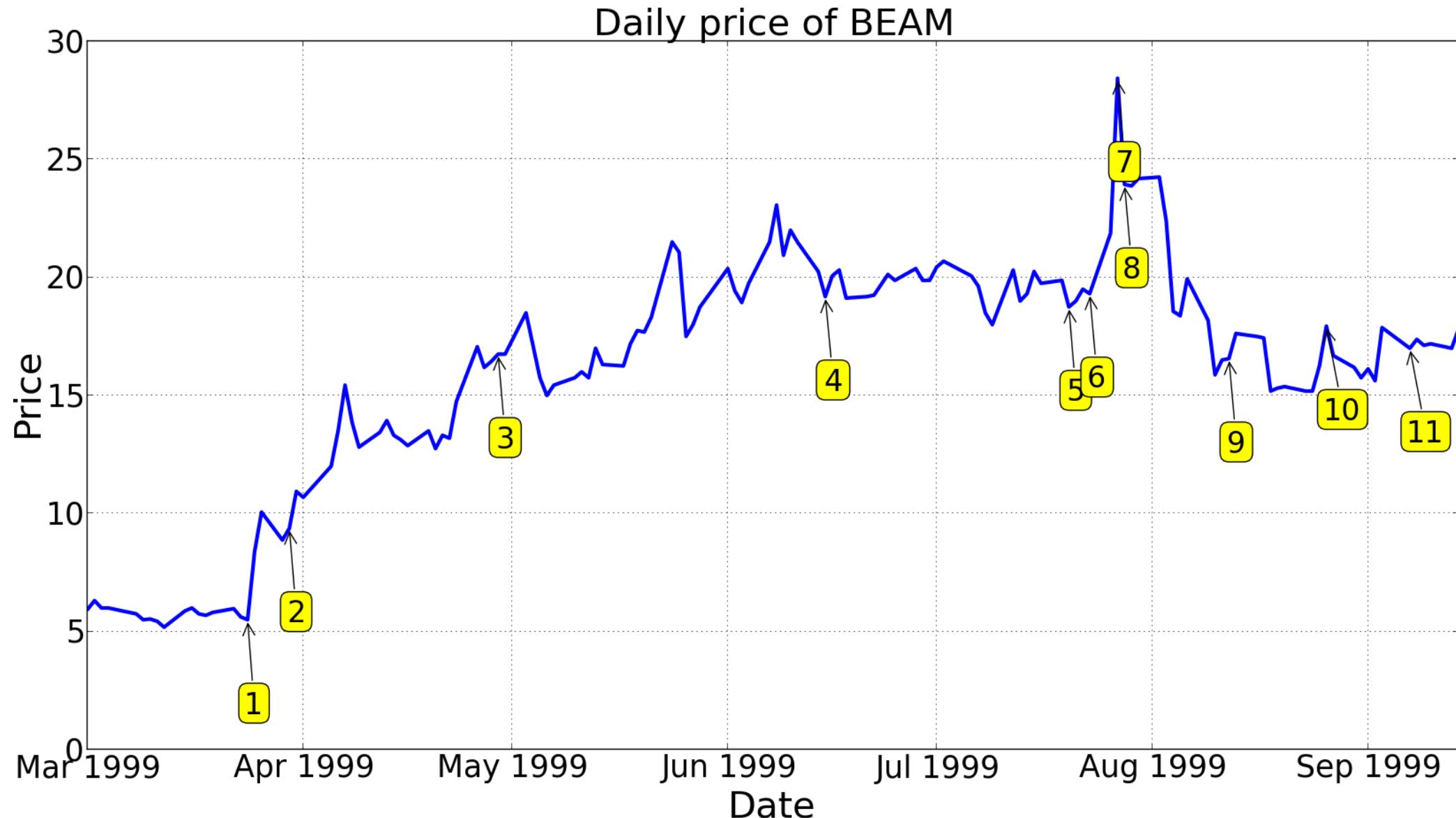
- **Task:** predict the stock market based on the stories that appear on the news wires

# Mining News Stories to Predict Stock Price Movement

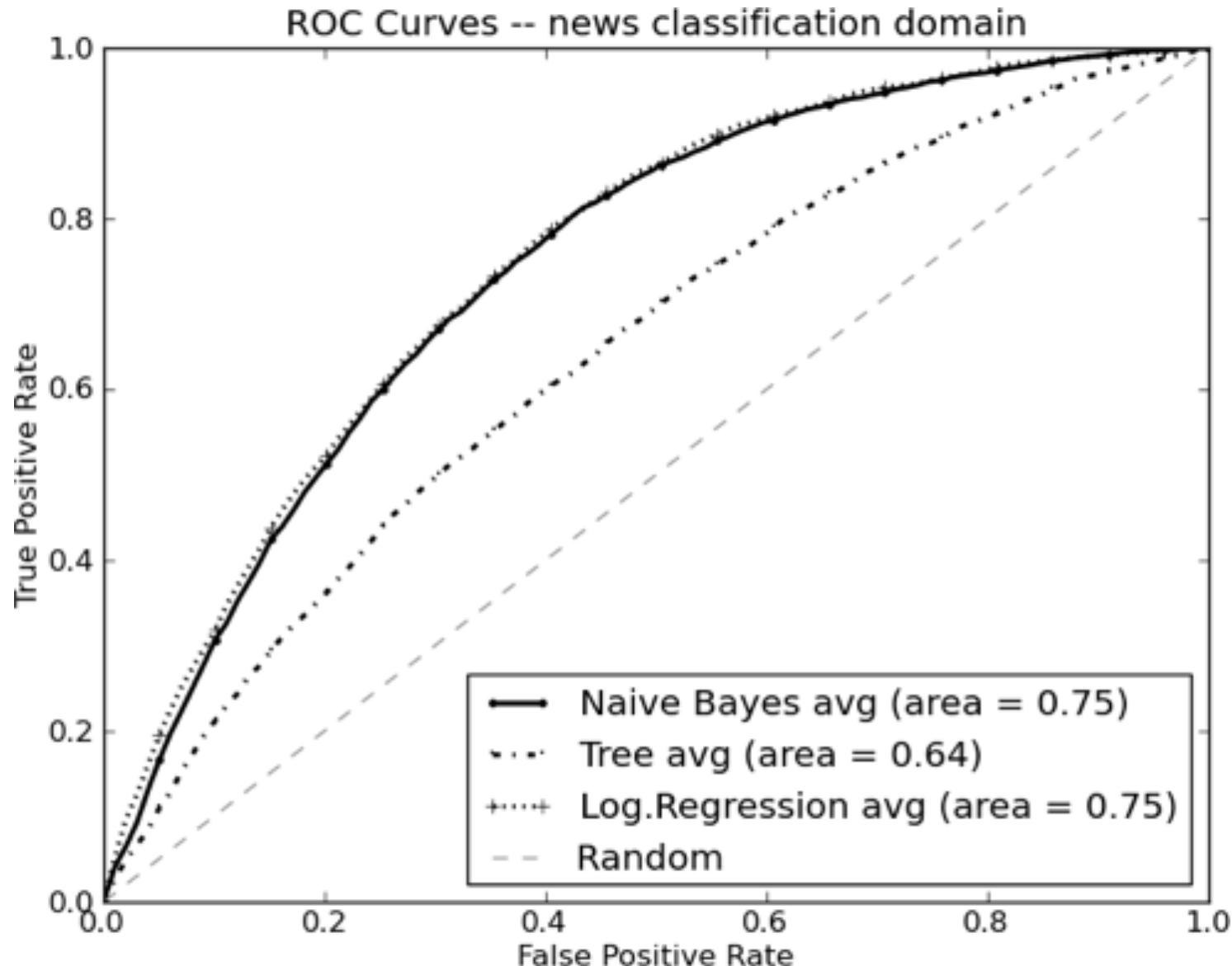
---



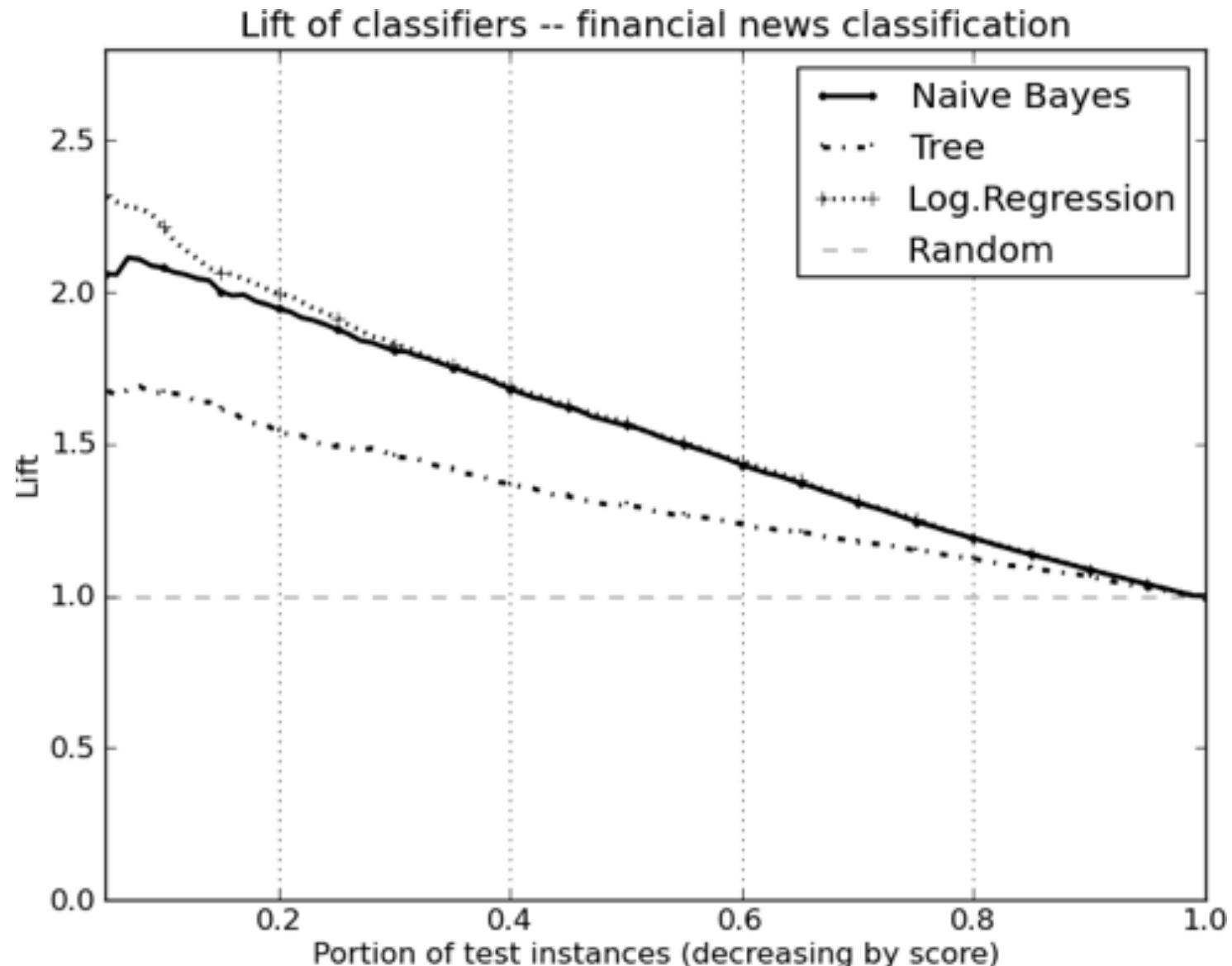
# Mining News Stories to Predict Stock Price Movement



# Mining News Stories to Predict Stock Price Movement



# Mining News Stories to Predict Stock Price Movement



# Mining News Stories to Predict Stock Price Movement

---

alert(s,ed), architecture, auction(s,ed,ing,eers), average(s,d), award(s,ed), bond(s), brokerage, climb(ed,s,ing), close(d,s), comment(ator,ed,ing,s), commerce(s), corporate, crack(s,ed,ing), cumulative, deal(s), dealing(s), deflect(ed,ing), delays, depart(s,ed), department(s), design(ers,ing), economy, econtent, edesign, eoperate, esource, event(s), exchange(s), extens(ion,ive), facilit(y,ies), gain(ed,s,ing), higher, hit(s), imbalance(s), index, issue(s,d), late(ly), law(s,ful), lead(s,ing), legal(ity,ly), lose, majority, merg(ing,ed,es), move(s,d), online, outperform(s,ance,ed), partner(s), payments, percent, pharmaceutical(s), price(d), primary, recover(ed,s), redirect(ed,ion), stakeholder(s), stock(s), violat(ing,ion,ors)