# Text Models

# Previously on CSCI 4622
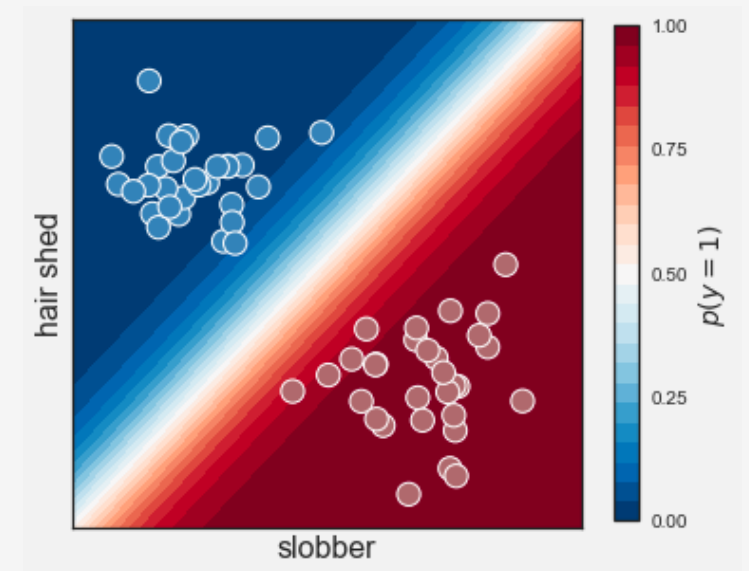
**Logistic Regression**:

○ Learn parameters $\beta_0, \beta_1, \ldots, \beta_p$ to model probability that example is in class by

$$p(y = 1 \mid \mathbf{x}) = \text{sigm}\left(\beta_0 + \beta_1 x_1 + \cdots + \beta_p x_p\right)$$

○ Decision rule for two-feature binary classification:

$$\hat{y} = \begin{cases} 1 & \text{if sigm}(\hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2) \geq 0.5 \\ 0 & \text{if sigm}(\hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2) < 0.5 \end{cases}$$

# Logistic Regression for Spam vs Ham

**Example**: How would you classify each of the following emails?

o **Email 1**: Mom, I got a job in Nigeria.

HAMMY

o **Email 2**: Jobs in Nigeria!  Money, money, money, money!

SPAMMY

⇒ lowercase

⇒ remove 's

# Vector Space Models of Text

Before we can use Logistic Regression, have to define what the features are.

$$p(y = 1 \mid \mathbf{x}) = \text{sigm}\left(\beta_0 + \beta_1 x_1 + \cdots + \beta_p x_p\right)$$

Most text models are what we call **vector space models**:

$$\text{"the quick brown fox"} \rightarrow \begin{bmatrix} 0 \\ 1 \\ 0 \\ 1 \\ 1 \\ 0 \end{bmatrix}$$

← QUICK

← BROWN

← FOX

# Vector Space Models of Text

Before we can use Logistic Regression, have to define what the features are.

$$p(y = 1 \mid \mathbf{x}) = \operatorname{sigm}\left(\beta_0 + \beta_1 x_1 + \cdots + \beta_p x_p\right)$$

Most text models are what we call **vector space models**:

$$\text{"the quick brown fox"} \rightarrow \begin{bmatrix} 0 \\ 1 \\ 0 \\ 1 \\ 1 \\ 0 \end{bmatrix} \qquad \text{"the fox jumps over the log"} \rightarrow \begin{bmatrix} 1 \\ 0 \\ 1 \\ 0 \\ 1 \\ 1 \end{bmatrix} \begin{matrix} \leftarrow \text{jumps} \\ \\ \leftarrow \text{over} \\ \\ \leftarrow \text{Fox} \\ \leftarrow \text{log} \end{matrix}$$

# Vector Space Models of Text

o  Suppose we have a training set comprised of many documents (articles, tweets, reviews)

o  From the training set we extract a vocabulary V of distinct words

o  Each document is represented by a feature vector $\mathbf{x}$ of length $p = |V|$

o  Each feature $x_k$ corresponds to a particular word in the vocabulary

**Example**: For the quick brown fox example, our vocabulary might be:

$$V = \{ \text{ quick, brown, fox, jump, over, log} \}$$

$|v| = 6$

1    2    3    4    5    6

# Vector Space Models of Text

**Example**: For the quick brown fox example, our vocabulary might be:

$$V = \{ \text{ quick}, \text{ brown}, \text{ fox}, \text{ jump}, \text{ over}, \text{ log}\}$$

o Represent a mapping from vocabulary to feature indices by a hash table

$$V = \{\text{quick} : 2, \text{ brown} : 4, \text{ fox} : 5, \text{ jump} : 1, \text{ over} : 3, \text{ log} : 6\}$$

**Binary Text Models** encode a document as a binary feature vector, where feature $x_k$ is 1 if term k appears anywhere in the document, and 0 otherwise.

# Vector Space Models of Text

The mapping for our quick brown fox example might be:

$$V = \{\text{quick} : 2, \ \text{brown} : 4, \ \text{fox} : 5, \ \text{jump} : 1, \ \text{over} : 3, \ \text{log} : 6\}$$

**Binary Text Models** encode a document as a binary feature vector, where feature $x_k$ is 1 if term k appears anywhere in the document, and 0 otherwise.

**Example**: Encode the document "That quick brown fox is quick!" using the binary model

$$x = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 1 \\ 1 \\ 0 \end{bmatrix} \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \\ 6 \end{matrix}$$

# Logistic Regression for Text

Suppose we've encoded a **training set** of documents and labels as $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ where the classes are represented as the labels $y_i \in \{0, 1\}$
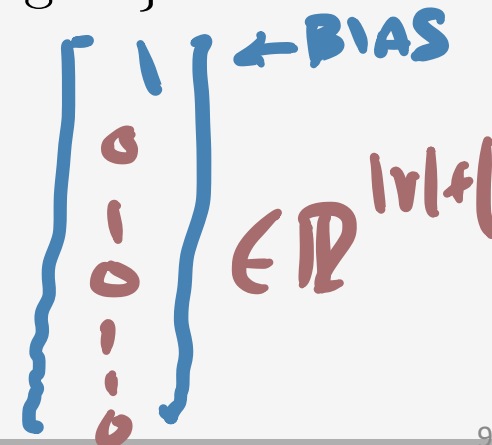
We want to learn a Logistic Regression model of the form

$$p(y = 1 \mid \mathbf{x}) = \text{sigm}\left(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_p x_p\right) = \text{sigm}(\boldsymbol{\beta}^T \mathbf{x})$$

Need to make vector lengths work out, so prepend each feature vector $\mathbf{x}$ with a 1

$$V = \{\text{bias} : 0, \ \text{quick} : 2, \ \text{brown} : 4, \ \text{fox} : 5, \ \text{jump} : 1, \ \text{over} : 3, \ \text{log} : 6\}$$

**Example**: The document "That quick brown fox is quick!" then becomes:

# Logistic Regression for Spam vs Ham

$sigm(\beta^T x)$

| feature | bias | $viagra$ | $mom$ | $job$ | $nigeria$ | $money$ |
|---------|------|----------|-------|-------|-----------|---------|
| parameter | $\beta_0$ | $\beta_1$ | $\beta_2$ | $\beta_3$ | $\beta_4$ | $\beta_5$ |
| learned value | 0.1 | 3.0 | −2.0 | −1.0 | 2.0 | 0.5 |

$\frac{1}{1+e^{-0.9}}$

**Example**: How would you classify the following email using the binary text model?

o **Email 1**: Mom, I got a job in Nigeria.

$$x = \begin{bmatrix} 1 \\ 0 \\ 1 \\ 1 \\ 0 \end{bmatrix}$$

$$\beta^T = [\,.1 \quad 3 \quad -2 \quad -1 \quad 2 \quad .5\,] \begin{bmatrix} 1 \\ 0 \\ 1 \\ 1 \\ 1 \\ 0 \end{bmatrix}$$

$$= .1 - 2 - 1 + 2 = -.9$$

$$P(y = SPAM \mid x) = sigm(-.9)$$

0.29

$< .5$

$\Rightarrow y = 0$

HAM

# Logistic Regression for Spam vs Ham

| feature | bias | *viagra* | *mom* | *job* | *nigeria* | *money* |
|---------|------|----------|-------|-------|-----------|---------|
| parameter | $\beta_0$ | $\beta_1$ | $\beta_2$ | $\beta_3$ | $\beta_4$ | $\beta_5$ |
| learned value | 0.1 | 3.0 | $-2.0$ | $-1.0$ | 2.0 | 0.5 |

*(handwritten above columns: 1, 0, 0, 1, 1, 1)*

**Example**: How would you classify the following email using the binary text model?

o **Email 2**: Jobs in Nigeria!  Money, money, money, money!

$$y = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$

$$\beta^T y = .1 - 1.0 + 2.0 + 0.5 = 1.6$$

$$P(y = SPAM \mid x) = sign(1.6) = 0.83$$

$$> .5 \Rightarrow \hat{y} = 1 \Rightarrow SPAM$$

# Logistic Regression for Spam vs Ham

| feature | bias | $viagra$ | $mom$ | $job$ | $nigeria$ | $money$ |
|---|---|---|---|---|---|---|
| parameter | $\beta_0$ | $\beta_1$ | $\beta_2$ | $\beta_3$ | $\beta_4$ | $\beta_5$ |
| learned value | 0.1 | 3.0 | $-2.0$ | $-1.0$ | 2.0 | 0.5 |

**Example**: How would you classify the following email using the binary text model?

o **Email 2**: Jobs in Nigeria!  Money, money, money, money!

Does something seem off about this?

The term "money" appeared 4 times in the document, but only gets a 1 in the feature vector

# Vector Space Models of Text

Recall our mapping for the quick brown fox example:

$$V = \{\text{bias} : 0, \text{jump} : 1, \text{quick} : 2, \text{over} : 3, \text{brown} : 4, \text{fox} : 5, \text{log} : 6\}$$

The **Bag-of-Words** Model takes into account the frequency of a term in a document

$$x_k = \# \text{ times term } k \text{ appears in document}$$

**Example**: Encode the document "That quick brown fox is quick!" using **Bag-of-Words**

$$\bar{x} = \begin{bmatrix} 1 \\ 0 \\ 2 \\ 0 \\ 1 \\ 1 \\ 0 \end{bmatrix}$$

# Logistic Regression for Spam vs Ham

*Email:* { } $\beta^T X = 0.1 \geqslant 0 \Rightarrow$ SPAM

| feature | bias | *viagra* | *mom* | *job* | *nigeria* | *money* |
|---|---|---|---|---|---|---|
| parameter | $\beta_0$ | $\beta_1$ | $\beta_2$ | $\beta_3$ | $\beta_4$ | $\beta_5$ |
| learned value | 0.1 | 3.0 | $-2.0$ | $-1.0$ | 2.0 | 0.5 |

**Example**: How would you classify the following email using the **Bag-of-Words** model?

o **Email 2**: Jobs in Nigeria!  Money, money, money, money!

$$y = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \\ 1 \\ 4 \end{bmatrix}$$

$\beta^T X = .1 - 1.0 + 2.0 + 4 \times .5$

$= 3.1$   $\text{Sigm}(3.1) = 0.98$

# Practical Implementation

Typically we store our training set in a matrix, where each row corresponds to a training example and each column corresponds to a feature.

Define the **Document-Term Matrix** $\mathbf{X}_{dt}$ for a **Bag-of-Words** model as follows:

$$[\mathbf{X}_{dt}]_{i,k} = \#\ \text{times term } k \text{ appears in document } i$$

**Example**: Suppose you have the following documents and vocabulary map, find $\mathbf{X}_{dt}$

```
Training Set :
d1 : new york new tribune
d2 : new york times
d3 : los angeles times
```

$$V = \{\text{new} : 3,\ \text{york} : 6,\ \text{tribune} : 5,\ \text{times} : 4,\ \text{los} : 2,\ \text{angeles} : 1\}$$

# Practical Implementation

**Example**: Suppose you have the following documents and vocabulary map, find $\mathbf{X}_{dt}$

Training Set :
d1 : new york new tribune
d2 : new york times
d3 : los angeles times

$V = \{\text{new} : 3,\ \text{york} : 6,\ \text{tribune} : 5,\ \text{times} : 4,\ \text{los} : 2,\ \text{angeles} : 1\}$

$$X_{dt} = \begin{array}{c} 1 \\ 2 \\ 3 \end{array} \begin{bmatrix} \overset{1}{0} & \overset{2}{0} & \overset{3}{2} & \overset{4}{0} & \overset{5}{1} & \overset{6}{1} \\ 0 & 0 & 1 & 1 & 0 & 1 \\ 1 & 1 & 0 & 1 & 0 & 0 \end{bmatrix}$$

# Practical Implementation

**Example**: Suppose you have the following documents and vocabulary map, find $\mathbf{X}_{dt}$

```
Training Set :
d1 : new york new tribune
d2 : new york times
d3 : los angeles times
```

$$V = \{\text{new} : 3, \ \text{york} : 6, \ \text{tribune} : 5, \ \text{times} : 4, \ \text{los} : 2, \ \text{angeles} : 1\}$$

**Solution**: We found

$$\mathbf{X}_{df} = \begin{bmatrix} 0 & 0 & 2 & 0 & 1 & 1 \\ 0 & 0 & 1 & 1 & 0 & 1 \\ 1 & 1 & 0 & 1 & 0 & 0 \end{bmatrix}$$

SPARSE MATRIX

**Question**: Suppose our vocabulary is huge, what data structure should we use for $\mathbf{X}_{dt}$ ?

# Practical Implementation

**Question**: What if you're have documents of widely varying lengths?  Should these be treated differently?

**Idea**: Rescale feature vectors so that longer docs don't overpower shorter docs

o  Scale rows of $\mathbf{X}_{dt}$  to have unit length

$$\mathbf{X}_{df} = \begin{bmatrix} 0 & 0 & 2 & 0 & 1 & 1 \\ 0 & 0 & 1 & 1 & 0 & 1 \\ 1 & 1 & 0 & 1 & 0 & 0 \end{bmatrix} \Rightarrow \mathbf{X}_{ndf} = \begin{bmatrix} 0 & 0 & 0.82 & 0 & 0.41 & 0.41 \\ 0 & 0 & 0.58 & 0.58 & 0 & 0.58 \\ 0.58 & 0.58 & 0 & 0.58 & 0 & 0 \end{bmatrix}$$
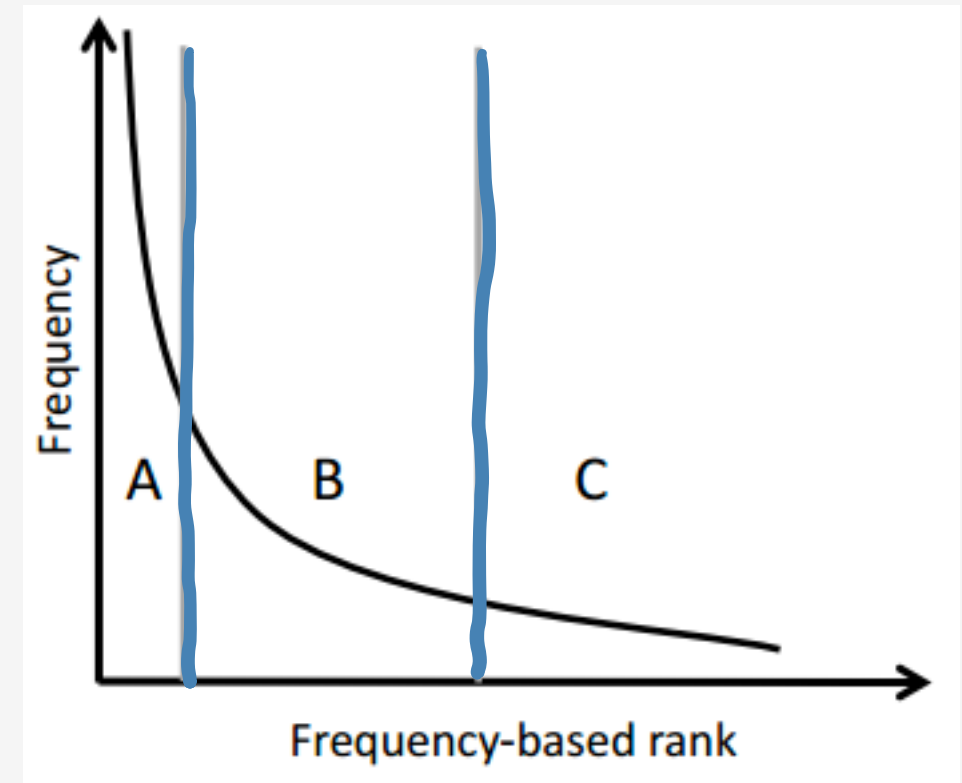
# Building Better Features

**Question**: What words do you think should be good for doing classification?

STOP WORDS = { is, the, that }

⇒ COMMON ⇒ LOW INFORMATION

# Building Better Features

**Question**: What words do you think should be good for doing classification?

o Column A: Words that are too common **can't discriminate** between classes.

o Column C: Words that are too uncommon **can't generalize** to new data

o Column B: Good features are words that are common, but not too common.

# Term Freq. - Inverse Document Freq.

**Idea**: Build term features based on how frequent a term is in a particular document Bnt and how many documents that term occurs in.

$$\texttt{tfidf(d,t))} = \texttt{tf(d,t)} \times \texttt{idf(t)}$$

# Term Freq. – Inverse Document Freq.

**Idea**: Build term features based on how frequent a term is in a particular document Bnt and how many documents that term occurs in.

$$\texttt{tfidf(d,t))} = \texttt{tf(d,t)} \times \texttt{idf(t)}$$

o The term frequency is the frequency of the term in the particular document

$$\texttt{tf(d,t)} = \texttt{\# times term t appears in document d}$$

# Term Freq. – Inverse Document Freq.

**Idea**: Build term features based on how frequent a term is in a particular document Bnt and how many documents that term occurs in.

$$\texttt{tfidf(d,t))} = \texttt{tf(d,t)} \times \texttt{idf(t)}$$

o The term frequency is the frequency of the term in the particular document

$$\texttt{tf(d,t)} = \texttt{\# times term t appears in document d}$$

o The inverse document frequency is measure of how many documents the term appears in

$$\texttt{idf(t)} = \log\left(\frac{1 + n_d}{1 + \texttt{df(} \, \texttt{t)}}\right) + 1$$

# Term Freq. - Inverse Document Freq.

The inverse document frequency is measure of how many documents the term appears in

o $n_d$ is the total number of training documents

o df(d,t) is the number of documents that contain at least one instance of term t

$$idf(t) = \log\left(\frac{1 + n_d}{1 + df(\blacksquare\ t)}\right) + 1$$

# Term Freq. – Inverse Document Freq.

**Idea**: Build term features based on how frequent a term is in a particular document Bent and how many documents that term occurs in.

$$\mathtt{tfidf(d,t)) = tf(d,t) \times idf(t)}$$

# Term Freq. – Inverse Document Freq.

**Example**: Compute the tfidf score for the word "new" in the following document set

```
Training Set:
d1:new york new tribune
d2:new york times
d3:los angeles times
```

WE'LL WORK THIS EXAMPLE OUT IN THE
LOGISTIC REGRESSION HANDS ON NOTEBOOK.

# Term Freq. - Inverse Document Freq.

**Example**: Compute the tfidf score for the word "new" in the following document set

```
Training Set:
d1:new york new tribune
d2:new york times
d3:los angeles times
```

# Text Models Wrap-Up

o Vector Space Models allow us to represent documents and their terms as feature vectors

o Vector Space Models come in many flavors, from simple binary to TF-IDF.

o Much fancier text models exist, like Google's **Word2Vec** Model

**Next Time**:

o Learn to do text encoding in Scikit-Learn

o Explore text models and logistic regression for predicting sentiment in movie reviews