# MULTINOMIAL CLUSTERING AND TEXT MODELS

# TEXT DATA

## Setting

Data set: A huge set of text documents (e.g. all books in a library). The entire set of texts is called a **corpus**.

Can we learn models from text which describe natural language?

## Terminology

We have to distinguish occurences of words in a document and *distinct* words in the dictionary. We refer to words regarded as entries of the dictionary as **terms**.

# EXAMPLE 2: SIMPLE TEXT MODEL

## Data

Suppose our data is a text document. We are given a dictionary which contains all terms occurring in the document.

## Documents as vectors of counts

We represent the document as

$$\mathbf{H} = (H_1, \ldots, H_d) \qquad \text{where } H_j = \# \text{ occurences of term } j \text{ in document.}$$

Note:

- $d$ is the number of all terms (distinct words) in the dictionary i.e. $d$ is identical for all documents.

- $n = \sum_j H_j$ can change from document to document.

# EXAMPLE 2: SIMPLE TEXT MODEL

## Multinomial model

To define a simple probabilistic model of document generation, we can use a multinomial distribution $P(\mathbf{H}|\mathbf{t}, n)$. That means:

- ▶ Each word in the document is sampled independently of the other words.

- ▶ The probabilities of occurrence are

$$\Pr\{\text{ word } = \text{ term } j \} = t_j .$$

## Implicit assumption

The assumption implicit in this model is that the probability of observing a document is completely determined by how often each term occurs; the order of words does not matter. This is called the **bag-of-words assumption**.

# CONTEXT

## Task

Can we predict the next word in a text?

## Context

In language, the co-occurence and order of words is highly informative. This information is called the **context** of a word.

**Example:** The English language has over 200,000 words.

▶ If we choose any word at random, there are over 200,000 possibilities.

▶ If we want to choose the next word in

$$\text{There is an airplane in the } \underline{\quad}$$

the number of possibilities is *much* smaller.

## Significance for statistical methods

Context information is well-suited for machine learning: By parsing lots of text, we can record which words occur together and which do not.

The standard models based on this idea are called *n-gram models*.

## Bigram model

A bigram model represents the conditional distribution

$$\mathrm{Pr}(\text{word}|\text{previous word}) =: \mathrm{Pr}(w_l|w_{l-1}) \ ,$$

where $w_l$ is the $l$th word in a text.

## Representation by multinomial distributions

A bigram model is a *family* of $d$ multinomial distributions, one for each possible previous word.

## Estimation

For each term $k$, find all terms in the corpus which are preceeded by $k$ and record their number of occurences in a vector

$$\mathbf{H}_k = (H_{k1}, \ldots, H_{kd}) \qquad \text{where } H_{kj} = \text{ number of times term } j \text{ follows on term } k$$

Then compute the maximum likelihood estimate $\hat{\mathbf{t}}_k$ from the sample $\mathbf{H}_k$.
**Note:** Both $j$ and $k$ run through $\{1, \ldots, d\}$.

# *N*-GRAM MODELS

## Multinomial representation of bigram

The distributions in the bigram model are:

$$\Pr(\text{word} = j | \text{previous word} = k) = P(H_j = 1 | \hat{\mathbf{t}}_k, n = 1)$$

where $P$ is the multinomial distribution. The entire bigram model is the set

$$\{ P( \, . \, | \hat{\mathbf{t}}_k, n = 1) \, | \, k = 1, \ldots, d \}$$

## *N*-gram models

More generally, a model conditional on the $(N - 1)$ previous words

$$\Pr(w_l | w_{l-1}, \ldots, w_{l-(N-1)})$$

is called an *N*-**gram model** (with the predicted word, there are $N$ words in total).

## Unigrams

The special case $N = 1$ (no context information) is the simple multinomial word probability model which we discussed first. This model is also called a **unigram model**.

# LEARNING SHAKESPEARE (1)

## Unigram Model

To him swallowed confess hear both.
Which. Of save on trail for are ay device
and rote life have

Every enter now severally so, let

Hill he late speaks; or! a more to leg less
first you enter

Are where exeunt and sighs have rise
excellency took of.. Sleep knave we. near;
vile like

## Bigram Model

What means, sir. I confess she? then all
sorts, he is trim, captain.

Why dost stand forth thy canopy, forsooth;
he is this palpable hit the King Henry.
Live king. Follow.

What we, hath got so she that I rest and
sent to scold and nature bankrupt, nor the
first gentleman?

Enter Menenius, if it so many good
direction found'st thou art a strong upon
command of fear not a liberal largess
given away, Falstaff! Exeunt

**Trigram Model**

Sweet prince, Falstaff shall die. Harry of Monmouth's grave.

This shall forbid it should be branded, if renown made it empty.

Indeed the duke; and had a very good friend.

Fly, and will rid me these news of price. Therefore the sadness of parting, as they say, 'tis done.

**Quadrigram Model**

King Henry. What! I will go seek the traitor Gloucester. Exeunt some of the watch. A great banquet serv'd in;

Will you not tell me who I am?

It cannot be but so.

Indeed the short and the long. Marry, 'tis a noble Lepidus.

# COMPLEXITY OF $N$-GRAM MODELS

## Enumerating contexts

An $N$-gram model considers ordered combinations of $N$ terms (=*distinct* words). Say a corpus contains 100,000 words. Then there are

$$100000^N = 10^{5N}$$

possible combinations.

## Naive estimate

If we require on average $n$ observations per combination to get a reliable estimate, we would need a corpus containing $n \cdot 10^{5N}$ words.

## Consequence

In practice, you typically encountner bigrams or trigrams. Research labs at some internet companies have reported results for higher orders.

## Task

Suppose we have a corpus consisting of two types of text, (1) cheap romantic novels and (2) books on theoretical physics. Can a clustering algorithm with two clusters automatically sort the books according to the two types?

(We will see that there is more to this than solving artificial sorting problems.)

## Clustering model

We assume the corpus is generated by a multinomial mixture model of the form

$$\pi(\mathbf{H}) = \sum_{k=1}^{K} c_k P(\mathbf{H}|\mathbf{t}_k) \; ,$$

i.e. each component $P(\mathbf{H}|\mathbf{t}_k)$ is multionmial.
**However:** We are now considering **documents** rather than individual words.

## Estimation

Apply EM algorithm for multinomial mixture models.

## Thought experiment

Say we run a mixture of two multinomial distributions on the cheap romantic novels and theoretical physics textbooks.

Outcome:

▶ Each cluster will roughly represent one of the two topics.

▶ The two parameter vectors $\mathbf{t}_1$ and $\mathbf{t}_2$ represent distributions of words in *texts of the respective topic*.

## Word distributions as topics

This motivates the interpretation of clusters as topics.

$$\mathbf{t}_k = \text{ distribution of words that characterizes topic } k$$

Language models derived from this idea are called **topic models**.