

# Distributed representations

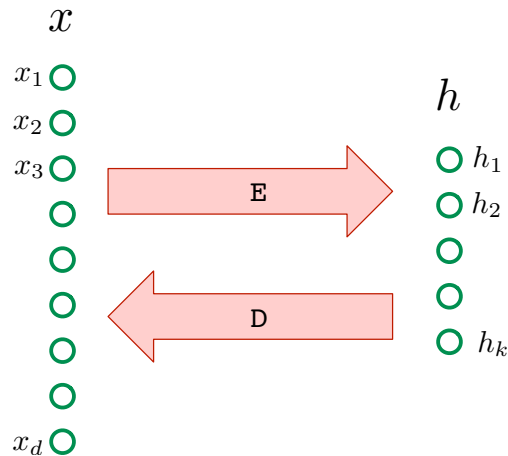
Sanjoy Dasgupta

University of California, San Diego

## Topics we'll cover

- ① One-hot versus distributed encodings
- ② Word embeddings

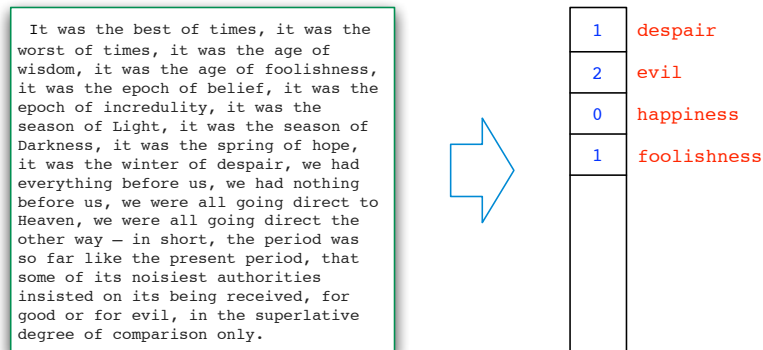
# One-hot versus distributed representations



- **k-means**: **one-hot** encoding
- **PCA**: **distributed** encoding

## The bag-of-words representation

One-hot encoding of words:



- Fix  $V$  = some vocabulary.
- Treat each sentence (or document) as a vector of length  $|V|$ :

$$x = (x_1, x_2, \dots, x_{|V|}),$$

where  $x_i = \#$  of times the  $i$ th word appears in the sentence.

## Word co-occurrences

*You shall know a word by the company it keeps.* (J.R. Firth, 1957)

- Much of the meaning of a word  $w$  is captured by the words it co-occurs with:

$w_1 \ w_2 \ w_3 \ w \ w_4 \ w_5 \ w_6$

- Find an embedding of words based on these co-occurrences.

## A simple approach to word embedding

Fix a vocabulary  $V$ . Then, using a corpus of text:

- ① Look at each word  $w$  and its surrounding *context*:  $w_1 \ w_2 \ w_3 \ w \ w_4 \ w_5 \ w_6$ 
  - $n(w, c) = \#$  times word  $c$  occurs in the context of word  $w$
  - Yields a probability distribution  $\Pr(c|w)$ .

- ② Positive pointwise mutual information:

$$\Phi_c(w) = \max \left( 0, \log \frac{\Pr(c|w)}{\Pr(c)} \right)$$

This is a  $|V|$ -dimensional representation of word  $w$ .

- ③ Reduce dimension using PCA.

## The embedding

- Which word's vector is closest to that of Africa?

Asia

- Solving analogy problems: king is to queen as man is to ?

- $\text{vec}(\text{king}) - \text{vec}(\text{queen}) = \text{vec}(\text{man}) - \text{vec}(?)$
- $\text{vec}(?) = \text{vec}(\text{man}) - \text{vec}(\text{king}) + \text{vec}(\text{queen})$
- Nearest neighbor of this vector is  $\text{vec}(\text{woman})$ .