

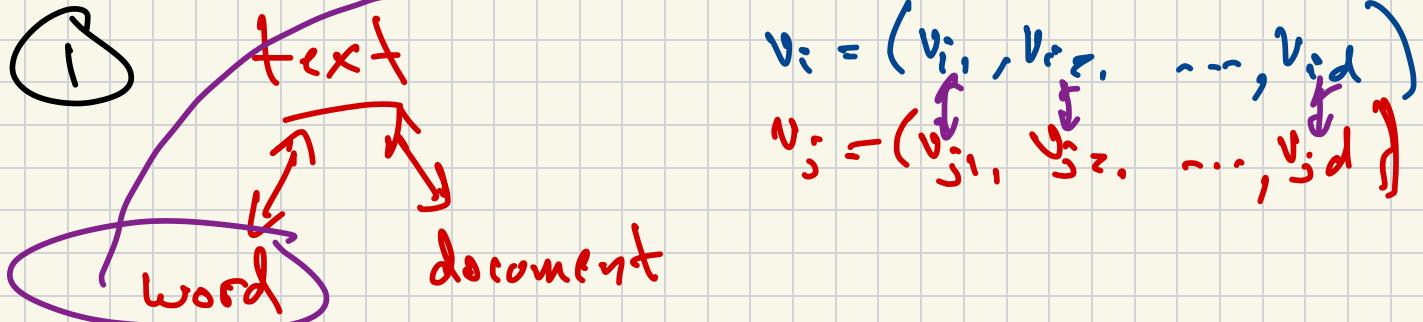
L3: Word Embeddings

Aug 25, 2025
Data Mining

Jeff Phillips

Data X , raw data $x_i \in X$

raw data \rightarrow vector $v_i \in \mathbb{R}^d$



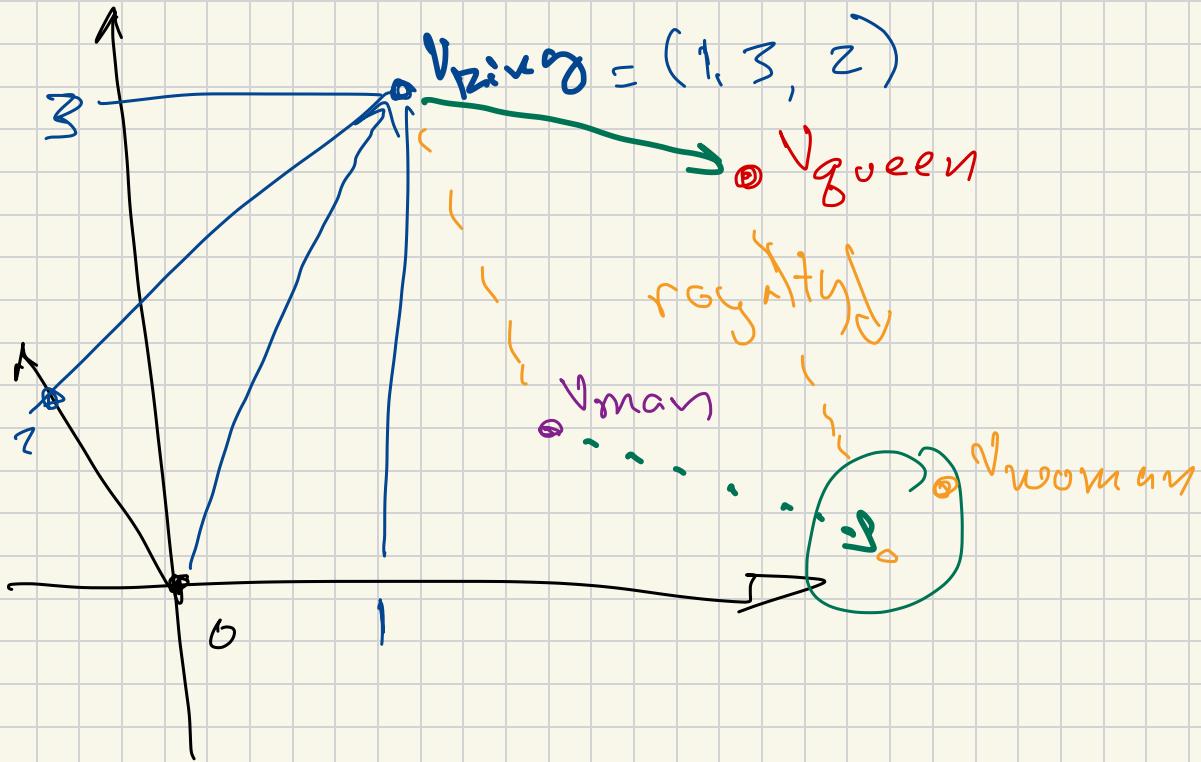
2 Distributional Hypothesis

alts: alphabetic, dewey-decimal, wordnet

\hookrightarrow words are similar if they tend to have similar context.

Word Vectors

$v_j, v_{king} \in \mathbb{R}^{d=300}$
 $j \in [n]$



M = 100,000 words

word = "octopus"

Corpus of text w, N words.

"I saw octopus"
"the larger octopus"

Teacher this weekend
she the fish at the
zoo.

co-occurrence vector $b_i \in \mathbb{R}^m$

$b_{j,i} = \# \text{ times word } j \text{ has word } i$
in its context.

Probability $P(j) = \frac{\# \text{ word } j \text{ occurs}}{N}$
 $P(i,j) = b_{i,j}/N$

co-occurrence vector $b_j \in \mathbb{R}^m$

$b_{j,i}$ = # times word j has word i
in its context.

probability $P(j) = \frac{\# \text{word } j \text{ occurs}}{N}$
 $P(i,j) = b_{j,i}/N$

pointwise mutual information $f_{i,j}$

$$\log \left(\frac{P(i,j)}{P(i) \cdot P(j)} \right)$$

$v_{j,i} = \max \left\{ 0, \log \left(\frac{P(j,i)}{P(i) \cdot P(j)} \right) \right\}$

vector $v_j = (v_{j,1}, v_{j,2}, \dots, v_{j,m}) \in \mathbb{R}^m$

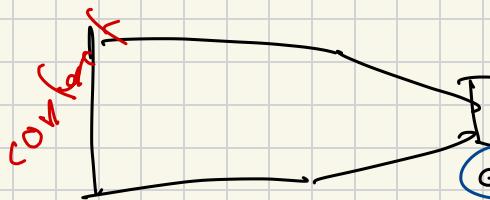
similar to tf-idf
BM25

Self-Supervised Learning (X, y)

The fish swim away from the shore. Next word predictions come from expert.

Mask out a word and predict its value.

Mask out a word
use context to predict its value.



$$\langle a | v_j \rangle$$

2013, 2014

Word2vec GLoVe

= each word has 1 representation

"bank"

↳ 2019 Elmo (LSTM)
If ϕ [Context + word] $\rightarrow v \in \mathbb{R}^d$

↳ Best (Transformer, attention)
context text