

# *Neural Language Model, RNNs*

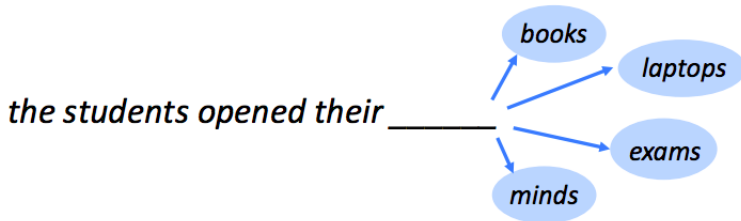
Pawan Goyal

CSE, IIT Kharagpur

CS60010

# Language Modeling

Language Modeling is the task of predicting what word comes next.



- **Goal:** Compute the probability of a sentence or sequence of words:

$$P(W) = P(w_1, w_2, w_3, \dots, w_n)$$

- **Related Task:** probability of an upcoming word:

$$P(w_4 | w_1, w_2, w_3)$$

- A model that computes either of these is called a **language model**

*Program LMs*

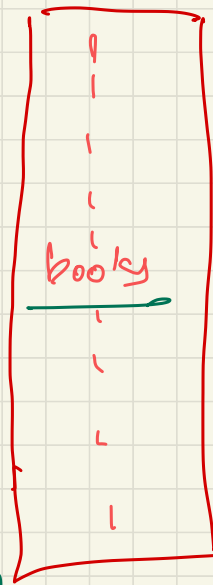
# Why should we care about language modeling?

- Language Modeling is a benchmark task that helps us measure our progress on understanding language
- Language Modeling is fundamental to many NLP tasks, especially those involving generating text or estimating the probability of text:
  - ▶ Predictive typing ✓
  - ▶ Speech recognition ✓
  - ▶ Handwriting recognition
  - ▶ Spelling/grammar correction

How'll I build an LM?

Large corpus

The students opened their



Vocabulary

[Every unique word  
in the corpus]

$P(\text{books} \mid \text{the students opened their})$

$$= \frac{\text{Count (the students opened their books)}}{\text{Count (the students opened their)}}$$

n-gram LM  $O(v) \leftarrow P(\text{books}) \rightarrow \text{unigram LM}$

$O(V^2) \in P(\text{bootstrap thets}) \rightarrow \text{bigrum LM}$   
 $O(V^3) \in P(\text{bootstrap opened thets}) \rightarrow \text{trigram}$

— storage issues

— Does not generalize

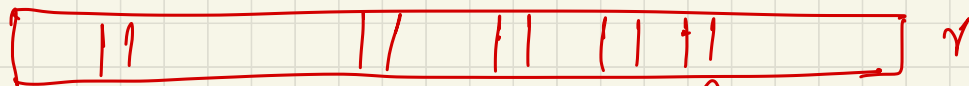
o/p layers ( $V$ -dim)

hidden ( $n$ )

o/p (the students opened their)

→ How to feed this as i/p?

27



books



opened



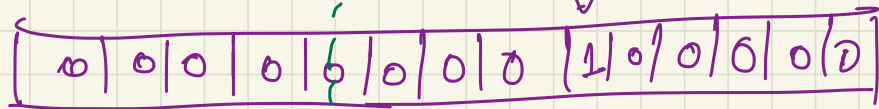
studied



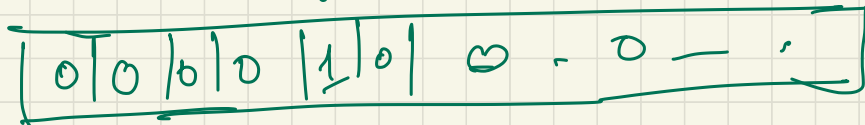
the



news



$w$

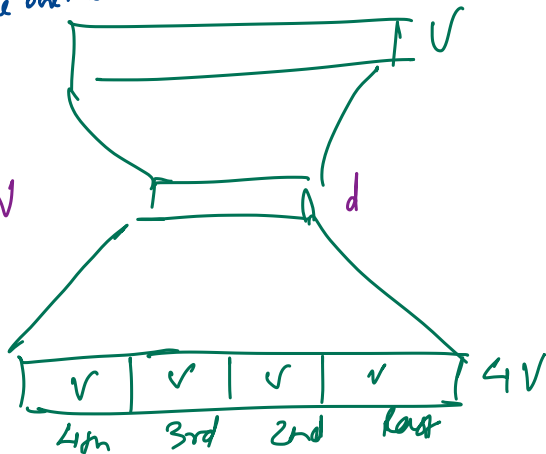


one-hot encoding  
r-dim

# A fixed-window neural language model

Students  
student  $\rightarrow$  separate one-hot vectors  
Can't generalize

$4V \times d + dV$   
 $= 5dV$   
exponential  $\rightarrow$  linear



~~as the proctor started the clock~~  
discard

the students opened their

fixed window  
pupils

# A fixed-window neural language model

output distribution

$$\hat{y} = \text{softmax}(Uh + b_2) \in \mathbb{R}^{|V|}$$

hidden layer

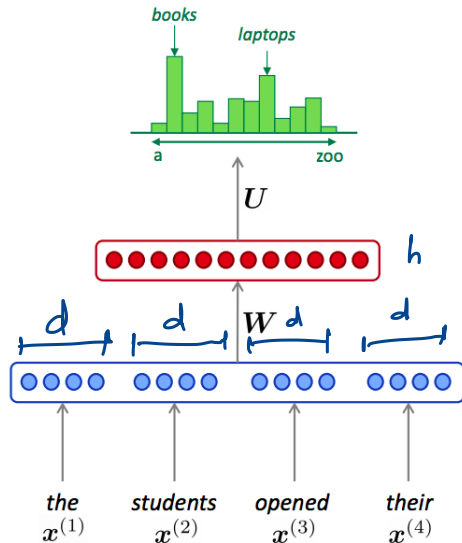
$$h = f(We + b_1)$$

concatenated word embeddings

$$e = [e^{(1)}; e^{(2)}; e^{(3)}; e^{(4)}]$$

words / one-hot vectors

$$x^{(1)}, x^{(2)}, x^{(3)}, x^{(4)}$$





# How do we obtain word representations?

In traditional NLP / IR, words are treated as discrete symbols.

## One-hot representation

Words are represented as one-hot vectors: one 1, the rest 0s

motel [0 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0] AND  
hotel [0 0 0 0 0 0 0 1 0 0 0 0 0 0 0 0] = 0

## What is the problem?

- Vector dimension = number of words in vocabulary (e.g., 500,000)
- The vectors are orthogonal, and there is no natural notion of similarity between one-hot vectors!

motel [0.1 0.5 | 1 1 1] d  
hotel [0.2 0.7 | ] d

Self-supervision:-

Take a large corpus

Given a word, predict surroundings

# Word2Vec – A distributed representation


## *Distributional representation – word embedding?*

Any word  $w_i$  in the corpus is given a distributional representation by an embedding

$$w_i \in \mathbb{R}^d$$

i.e., a  $d$ –dimensional vector, which is mostly learnt!

*linguistics* = 
$$\begin{bmatrix} 0.286 \\ 0.792 \\ -0.177 \\ -0.107 \\ 0.109 \\ -0.542 \\ 0.349 \\ 0.271 \end{bmatrix}$$



# Learning Word Vectors: Overview

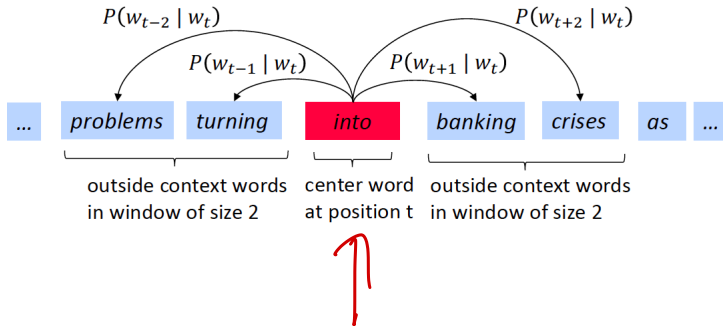
in-vector  $V \times d$  out-vector  $V \times d$

## Basic Idea: Use self-supervision

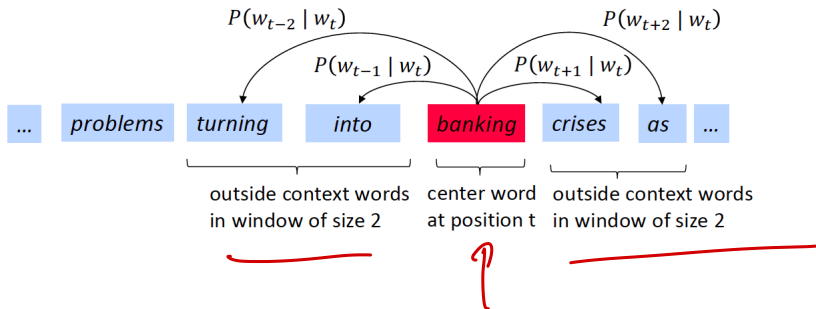
- We have a large corpus of text
- Every word in a fixed vocabulary is represented by a *vector*
- Go through each position  $t$  in the text, which has a center word  $c$  and context (“outside”) words  $o$
- Use the similarity of the word vectors for  $c$  and  $o$  to calculate the probability of  $o$  given  $c$  (or vice versa)
- Keep adjusting the word vectors to maximize this probability

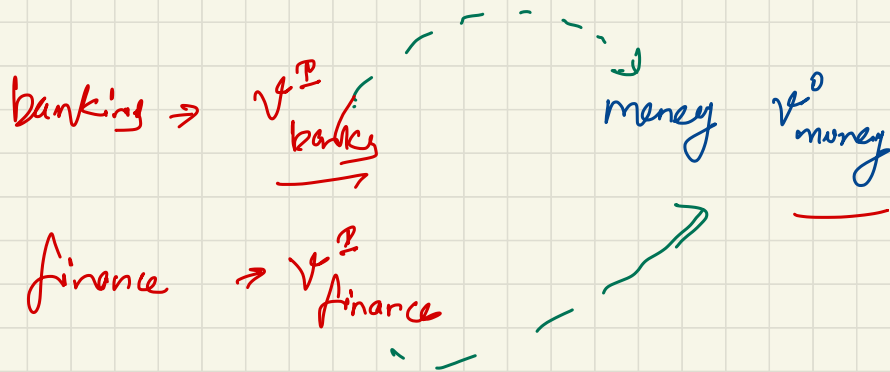
# Word2Vec (Skip-gram) Overview

Example windows and process for computing  $P(w_{t+j}|w_t)$



Example windows and process for computing  $P(w_{t+j}|w_t)$





hot

cold

