# *Neural Language Model, RNNs*

Pawan Goyal
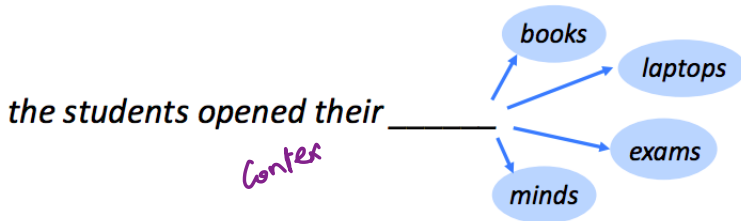
CSE, IIT Kharagpur

CS60010

Predictive typing

# Language Modeling

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

the students opened their _____

*Conter*

- books
- laptops
- exams
- minds

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

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

*chain rule*

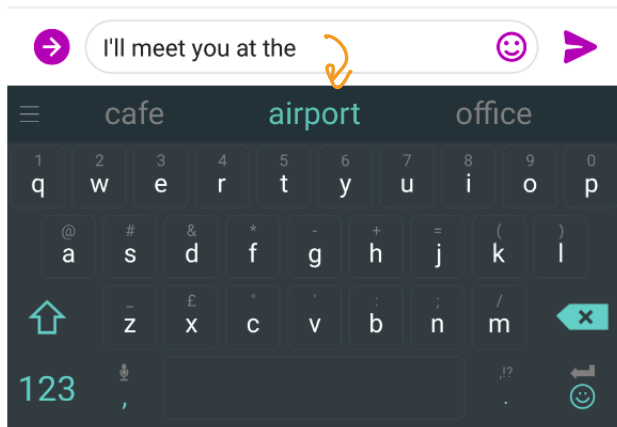- **Related Task:** probability of an upcoming word:

$$P(w_4 | w_1, w_2, w_3) \quad \to LM$$

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

# Language Modeling

- You can also think of a language model as a system that assigns probability to a piece of text.
- For example, if we have some text $x^{(1)}, \ldots, x^{(T)}$, then the probability of this text (according to the Language Model) is:

$$P(\boldsymbol{x}^{(1)}, \ldots, \boldsymbol{x}^{(T)}) = P(\boldsymbol{x}^{(1)}) \times P(\boldsymbol{x}^{(2)} | \boldsymbol{x}^{(1)}) \times \cdots \times P(\boldsymbol{x}^{(T)} | \boldsymbol{x}^{(T-1)}, \ldots, \boldsymbol{x}^{(1)})$$

$$= \prod_{t=1}^{T} P(\boldsymbol{x}^{(t)} | \boldsymbol{x}^{(t-1)}, \ldots, \boldsymbol{x}^{(1)})$$

This is what our LM provides

LM

# 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

*translation*

*chatbot*

*Lot of data → compress info → Parameter Set*

*the students opened their* _____ **books**

$Pr(books \mid their)$

$Pr(books \mid opened\ their)$

**Question**: How to learn a Language Model?

**Answer** (pre- Deep Learning): learn a *n-gram Language Model*!

**Definition:** A *n-gram* is a chunk of *n* consecutive words.

- unigrams: "the", "students", "opened", "their"
- bigrams: "the students", "students opened", "opened their"
- trigrams: "the students opened", "students opened their"
- 4-grams: "the students opened their"

**Idea:** Collect statistics about how frequent different n-grams are, and use these to predict next word.

$k = 101$     $V^{101}$     $100,000$

# n-gram language models

- First we make a simplifying assumption: $x^{(t+1)}$ depends only on the preceding *n-1* words.

*n-1* words

n-gram

$$P(x^{(t+1)} | x^{(t)}, \ldots, x^{(1)}) \cong P(x^{(t+1)} | x^{(t)}, \ldots, x^{(t-n+2)})$$ (assumption)

prob of a n-gram

prob of a (n-1)-gram

$$= \frac{P(x^{(t+1)}, x^{(t)}, \ldots, x^{(t-n+2)})}{P(x^{(t)}, \ldots, x^{(t-n+2)})}$$ (definition of conditional prob)

- **Question:** How do we get these *n*-gram and (*n*-1)-gram probabilities?
- **Answer:** By counting them in some large corpus of text!

$$\approx \frac{\text{count}(x^{(t+1)}, x^{(t)}, \ldots, x^{(t-n+2)})}{\text{count}(x^{(t)}, \ldots, x^{(t-n+2)})}$$ (statistical approximation)

Suppose we are learning a 4-gram Language Model.

~~as the proctor started the clock, the~~ *students opened their* ____ books

discard

condition on this

$$P(\boldsymbol{w}|\text{students opened their}) = \frac{\text{count}(\text{students opened their } \boldsymbol{w})}{\text{count}(\text{students opened their})}$$

For example, suppose that in the corpus:
- "students opened their" occurred 1000 times
- "students opened their books" occurred 400 times
  - → P(books | students opened their) = 0.4
- "students opened their exams" occurred 100 times
  - → P(exams | students opened their) = 0.1

Should we have discarded the "proctor" context?

**Storage**: Need to store count for all *n*-grams you saw in the corpus.

$$P(\boldsymbol{w}|\text{students opened their}) = \frac{\text{count(students opened their } \boldsymbol{w})}{\text{count(students opened their)}}$$

Increasing *n* or increasing corpus increases model size!

# A fixed-window neural language model

o/p    $w \in v$   [rectangle]   V

$gdh + vh$

[rectangle] h

I/p   [Students]   [opened]   [their]   3d   i/p size

1-hot    d-dim

$[3Vh + vh]$

~~as the proctor started the clock~~   the   students   opened   their   _____

discard      fixed window

# A fixed-window neural language model



**output distribution**
$$\hat{y} = \text{softmax}(\boldsymbol{U}\boldsymbol{h} + \boldsymbol{b}_2) \in \mathbb{R}^{|V|}$$

**hidden layer**
$$\boldsymbol{h} = f(\boldsymbol{W}\boldsymbol{e} + \boldsymbol{b}_1)$$

**concatenated word embeddings**
$$\boldsymbol{e} = [\boldsymbol{e}^{(1)}; \boldsymbol{e}^{(2)}; \boldsymbol{e}^{(3)}; \boldsymbol{e}^{(4)}]$$

**words / one-hot vectors**
$$\boldsymbol{x}^{(1)}, \boldsymbol{x}^{(2)}, \boldsymbol{x}^{(3)}, \boldsymbol{x}^{(4)}$$

# 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 1 0 0 0 0] AND
hotel [0 0 0 0 0 0 0 1 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!

# Word2Vec – A distributed representation

## Distributional representation – word embedding?

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

$$w_i \in R^d$$

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

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

dimensions?

If we label the dimensions in a hypothetical word vector (there are no such pre-assigned labels in the algorithm of course), it might look a bit like this:



| | King | Queen | Woman | Princess | ... |
|---|---|---|---|---|---|
| Royalty | 0.99 | 0.99 | 0.02 | 0.98 | |
| Masculinity | 0.99 | 0.05 | 0.01 | 0.02 | |
| Femininity | 0.05 | 0.93 | 0.999 | 0.94 | |
| Age | 0.7 | 0.6 | 0.5 | 0.1 | |

self-supervision

*Such a vector represents the 'meaning' of a word in some abstract way*

# Unsupervised

Test data

No labels

[no loss fn]

# Self-supervised

Train data

[No human-generated labels]

Generates labels from the data itself
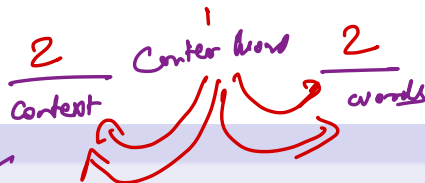
# Supervised

Training data with labels

*2 Context* — *Center Word* — *2 words*
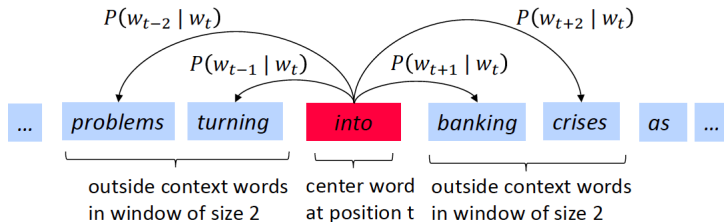
**Basic Idea: Use self-supervision**

- 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

\* word vectors are your params

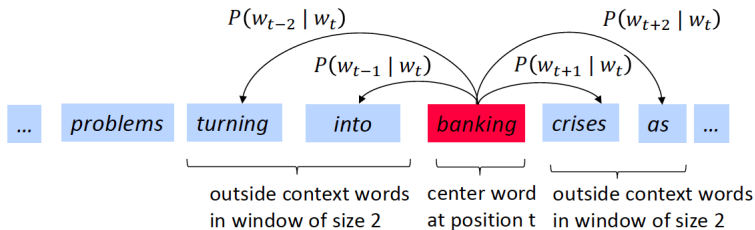Example windows and process for computing $P(w_{t+j}|w_t)$

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

# Word2Vec: objective function

We want to minimize the loss function:

$$J(\theta) = -\frac{1}{T} \sum_{t=1}^{T} \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \log P\left(w_{t+j} \mid w_t; \theta\right)$$

*2 v*     *vectors*

$V_w^g$    $U_w$

*How to calculate $P(w_{t+j}|w_t;\theta)$?*

We will use two vectors per word $w$:

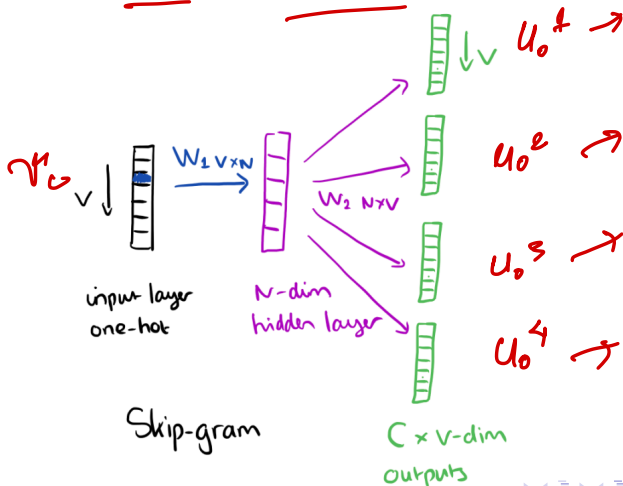- $v_w$ when $w$ is a center word
- $u_w$ when $w$ is a context word
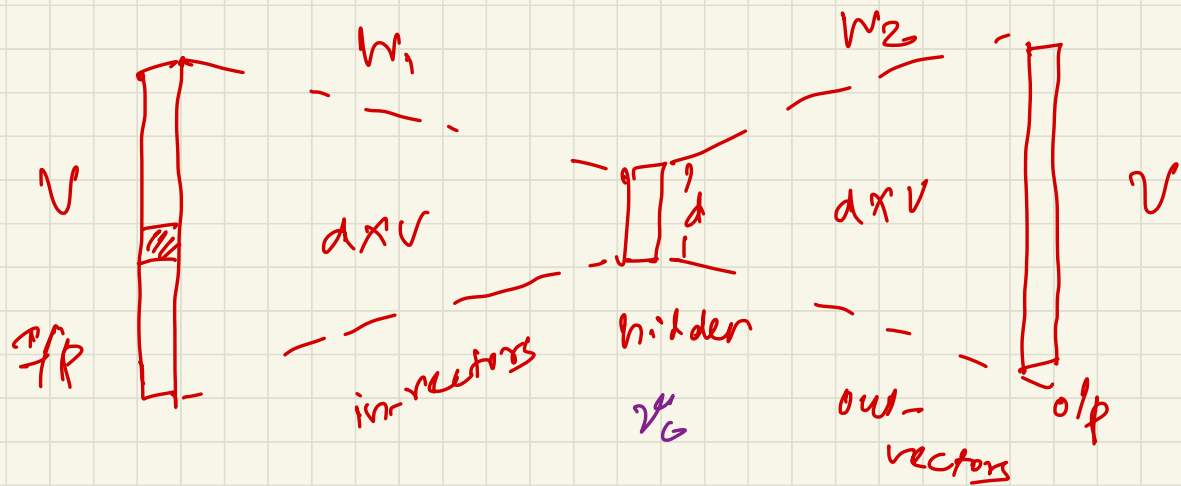
Then, for a center word $c$ and a context word $o$

$V_c^g$    $U_o$

$$P(o|c) = \frac{exp(u_o^T v_c)}{\sum_{w \in V} exp(u_w^T v_c)}$$

$$P(o|c) = \frac{exp(u_o^T v_c)}{\sum_{w \in V} exp(u_w^T v_c)}$$



input layer
one-hot

Skip-gram

N-dim
hidden layer

C × V-dim
outputs

$V$

$W_1$

$W_2$

$V$

$d \times V$

$d$

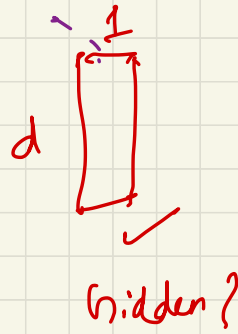$d \times V$

i/p

in-vectors

hidden

$V_G$

out-vectors

o/p
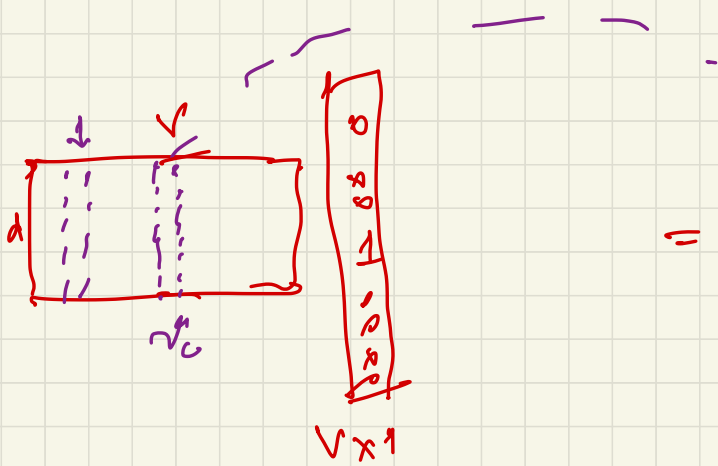
$\dfrac{C}{\text{bankity}}$

$\dfrac{O}{\text{Coises}}$

$W \rightarrow \dfrac{V_W^q + U_W}{d\text{-dim} \quad d\text{-dim}}$

$\boxed{2dV}$

$$V \times 1$$

$$=$$

$$d$$

hidden?

$$-\log p(o|c)$$

$$d \quad u_1^T$$

$$1$$
$$2$$
$$V$$
$$V$$

$$v_c$$

$$=$$

$$\frac{v_c}{}$$

$$\to u_1^T v_c$$

$$u_0^T v_c \xrightarrow{softmax} \exp(u_0^T v_c)$$

$$\underbrace{\sum \exp(u_w^T v_c)}_{w \in t}$$

### Skip-gram

Suppose you are computing the word vectors using Skip-gram architecture. You have 5 words in your vocabulary, $\{passed, through, relu, activation, function\}$ in that order and suppose you have the window, '*through relu activation*' in your corpora. You use this window with 'relu' as the center word and one word before and after the center word as your context.

### Compute the loss

Also, suppose that for each word, you have 2-dim in and out vectors, which have the same value at this point given by [1,-1],[1,1],[-2,1],[0,1],[1,0] for the 5 words, respectively. As per the Skip-gram architecture, the loss corresponding to the target word "activation" would be $-log(x)$. What is the value of $x$?

- Compute partial derivative of the loss with respect to $v_c$

$$\frac{\partial}{\partial v_c} \left( \quad \right)$$