Neural Language Model

Pawan Goyal

CSE, IIT Kharagpur

NLP - Interaction Hour

denied the

$$\boxed{0 \mid 0 \mid 0 \mid 0 \mid 0.1 \mid 0.2 \mid 0.3 \mid 0.4 \mid 0 \mid 0 \mid 0}$$

$V$

Smoothing

P(w | denied the)

$\boxed{1/|V| \quad \dots}$

unigram-prior smoothing

Voc $\rightarrow$ 1 million

500,000
words with
mynourfes.

Sanskrit SLP1

WX notation

$\frac{c}{=} \rightarrow t$

$\frac{c}{=} \rightarrow T$

There are names of some countries in South America, written in the Georgian language, together with their translations to English:

Brazilia

ბრაზილია — Brazil
პერუ — Peru
ურუგვაი — Uruguay

जापान — Japan
चीन — China
रूस — Russia

არგენტინა → Argentina
კოლუმბია → Columbia

What are the names, in English, of the two untranslated countries?

कमला — Sanskrit
Kamala

कमल

Hindi
↓
kamal

$i \rightarrow n$
$y$

कमल → Kamalo

Speech and Language Processing, 3rd Edition (draft): Chapter 9.
Sequence Processing with Recurrent Networks
https://web.stanford.edu/~jurafsky/slp3/9.pdf
*Many of these slides have been adapted from CS 224n. Some material is from the Deep Learning book.*

# Language Modeling

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

the students opened their _____ → books, laptops, exams, minds

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

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

# Language Modeling

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

the students opened their _____

books

laptops

exams

minds

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

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

*auto-completion*

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

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

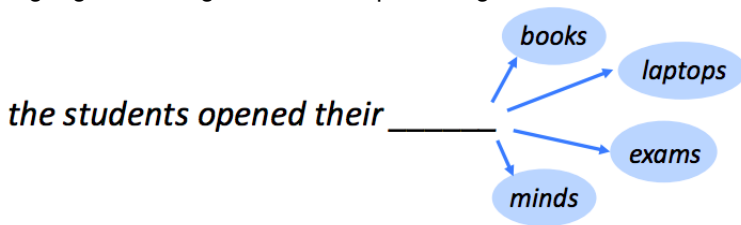## Language Modeling

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



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

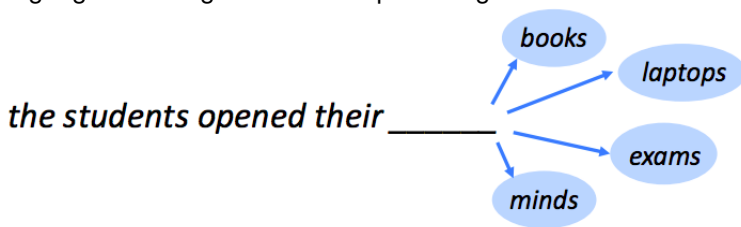$$P(W) = P(w_1, w_2, w_3, \ldots, w_n)$$

- **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**

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

# n-gram language models

*the students 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.

- uni grams: "the", "students", "opened", "their"
- bi grams: "the students", "students opened", "opened their"
- tri grams: "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.

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

$$P(x^{(t+1)}|x^{(t)},\ldots,x^{(1)}) = P(x^{(t+1)}|\underbrace{x^{(t)},\ldots,x^{(t-n+2)}}_{n\text{-}1 \text{ words}})$$

(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)

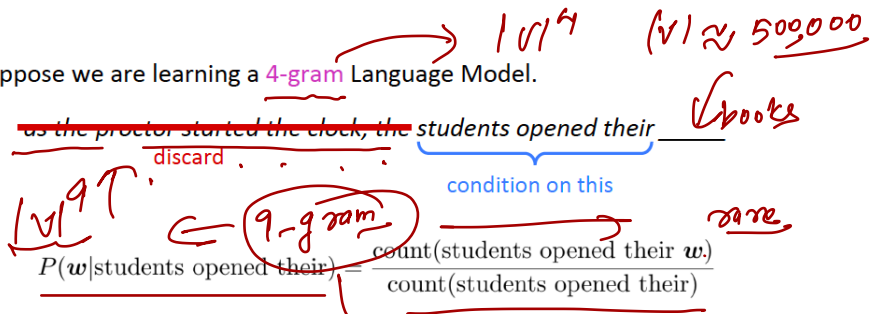- **Question:** How do we get these *n*-gram and (*n*-1)-gram probabilities?
- **Answer:** By counting them in some large corpus of text!

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

(statistical approximation)

$$B_i \qquad P(x_2|x_1) \qquad P'(x_1,x_2) \,/\, C(x_1)$$

Suppose we are learning a **4-gram** Language Model.

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

discard

condition on this

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

How many parameters

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.

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

Increasing *n* or increasing corpus increases model size!

2003
Bengio

C.E loss

o/p  $d$-dim $P_x$

$W$-dim

$U$

hidden layer

$h$

$$U^T\left(\sigma\left(Wx+b\right)\right)+c$$

I/player

$W_1$

$x$

$4d$

as  the  proctor  started  the  clock

discard

the   students   opened   their

fixed window

Adv. over ngram LM ?



books     laptops

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

$\boldsymbol{U}$     + hxV

$O(NI)$  $\Leftarrow$

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

$\boldsymbol{W}$ → Hd×h
Hd

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

Sachin  Tendulkor  will  play

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

the $\boldsymbol{x}^{(1)}$     students $\boldsymbol{x}^{(2)}$     opened $\boldsymbol{x}^{(3)}$     their $\boldsymbol{x}^{(4)}$

$play \uparrow \quad played \downarrow$

$d \times h \qquad d \times h \qquad d \times h \qquad d \times h$

$\uparrow \qquad \uparrow \qquad \uparrow$

$4d \times h$

$x^{(1)} \qquad x^{(2)} \qquad x^{(3)} \qquad x^{(4)}$

Sachin    Tendulkar   will

Sachin   Tendulkar   will   Sourly

# A fixed-window neural language model

**Improvements** over $n$-gram LM:
- No sparsity problem
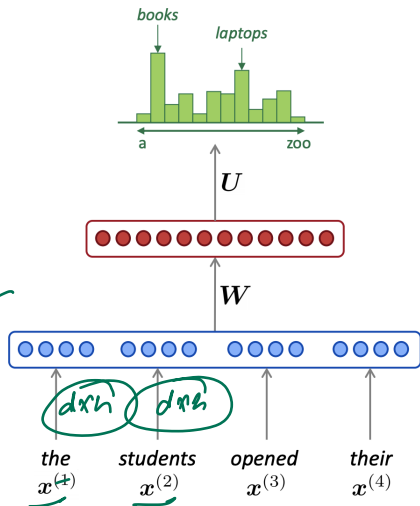- Don't need to store all observed $n$-grams

Remaining **problems**:
- Fixed window is too small
- Enlarging window enlarges $W$
- Window can never be large enough!
- $x^{(1)}$ and $x^{(2)}$ are multiplied by completely different weights in $W$. No symmetry in how the inputs are processed.

We need a neural architecture that can process *any length input*



$U$

$W$

*the*
$x^{(1)}$

*students*
$x^{(2)}$

*opened*
$x^{(3)}$

*their*
$x^{(4)}$

books

laptops

a                    zoo

# Recurrent Neural Networks



## Core Idea
Apply the same weights repeatedly!

We can process a sequence of vectors $x$ by applying a recurrence formula at each step:
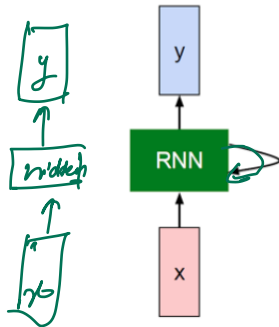
We can process a sequence of vectors $x$ by applying a recurrence formula at each step:

$$h_t = f_W(h_{t-1}, x_t)$$

new state — $h_t$

some function with parameters W — $f_W$

old state — $h_{t-1}$

input vector at some time step — $x_t$

Notice: the same function and the same set of parameters are used at every time step.

## Activation function for the hidden units

Assume the hyperbolic tangent activation function



$$tanh \left( U x^{(1)} + b \right) + c$$

$$y^{(1)}$$
$$\uparrow V^{\top}$$

$$h^{(1)} \xleftarrow{W} h^{(2)} \xrightarrow{W} h^{(3)}$$

$$\uparrow U \qquad \uparrow U$$

$$\frac{d}{x^{(1)}} \qquad \frac{d}{x^{(2)}} \qquad x^{(3)}$$

$$y^{(t)}$$
$$V^{\top} h^{(t)} + c$$

$$h^{(t)}$$

$$h^{(t)} = tanh \left( U x^{(t)} + W h^{(t-1)} + b \right)$$

$$- - - \quad x^{(t)}$$

$$h^{(t)} = \tanh\left(W h^{(t-1)} + U x^{(t)} + b\right)$$

$y^{(t)}$

$y^{(t)} \uparrow V$

softmax

$$y^{(t)} = \left(C + V h^{(t)}\right)$$

$y^{(t)}$

$h^{(t-1)} \longrightarrow \boxed{W} \rightarrow h^{(t)}$

$\uparrow$

$\boxed{h^{(t)}}$

$U$ :

$\dfrac{V}{J}, \dfrac{W}{L}, \dfrac{U}{l}$

$W : U$

$n \times |V| \quad n \times n \quad h \times d$

$\uparrow$

$\boxed{x^{(t)}}$

$\boxed{h^{(t-1)} \quad | \quad x^{(t)}}$

# Forward propagation for the RNN: first model

## Activation function for the hidden units

Assume the hyperbolic tangent activation function

## Form of output and loss function

Assume output is discrete - predicting words
We can obtain a vector normalized probabilities over the output - $\hat{y}$.

# *Forward propagation for the RNN: first model*

### *Activation function for the hidden units*

Assume the hyperbolic tangent activation function

### *Form of output and loss function*

Assume output is discrete - predicting words

We can obtain a vector normalized probabilities over the output - $\hat{y}$.

### *Update Equations*

Initial state - $h^{(0)}$

# Forward propagation for the RNN: first model

## Activation function for the hidden units

Assume the hyperbolic tangent activation function

## Form of output and loss function

Assume output is discrete - predicting words
We can obtain a vector normalized probabilities over the output - $\hat{y}$.

## Update Equations

Initial state - $h^{(0)}$
From $t = 1$ to $t = \tau$, the following update equation is applied:

$$a^{(t)} = b + Wh^{(t-1)} + Ux^{(t)}$$

$$a^{(t)} = b + Wh^{(t-1)} + Ux^{(t)}$$

$$h^{(t)} = tanh(a^{(t)})$$

$$o^{(t)} = c + Vh^{(t)}$$

$$\hat{y}^{(t)} = softmax(o^{(t)})$$

*pred.*

$$a^{(t)} = b + Wh^{(t-1)} + Ux^{(t)}$$

$$h^{(t)} = tanh(a^{(t)})$$

$$o^{(t)} = c + Vh^{(t)}$$

$$\hat{y}^{(t)} = softmax(o^{(t)})$$

This maps an input sequence to an output sequence of the same length.

Total loss is sum of the losses over all the time steps.
So, if $L^{(t)}$ is the negative log likelihood of $y^{(t)}$ given $x^{(1)}, \ldots, x^{(\tau)}$, then

Total loss is sum of the losses over all the time steps.

So, if $L^{(t)}$ is the negative log likelihood of $y^{(t)}$ given $x^{(1)}, \ldots, x^{(\tau)}$, then

$$
\begin{aligned}
& L(\{x^{(1)}, \ldots, x^{(\tau)}\}, \{y^{(1)}, \ldots, y^{(\tau)}\}) \\
=\ & \sum_t L^{(t)} \\
=\ & -\sum_t \log p_{model}(y^{(t)} | \{x^{(1)}, \ldots, x^{(\tau)}\})
\end{aligned}
$$

correct class

Total loss is sum of the losses over all the time steps.

So, if $L^{(t)}$ is the negative log likelihood of $y^{(t)}$ given $x^{(1)}, \ldots, x^{(\tau)}$, then

$$
\begin{aligned}
& L(\{x^{(1)}, \ldots, x^{(\tau)}\}, \{y^{(1)}, \ldots, y^{(\tau)}\}) \\
= \ & \sum_t L^{(t)} \\
= \ & -\sum_t \log p_{model}(y^{(t)} | \{x^{(1)}, \ldots, x^{(\tau)}\})
\end{aligned}
$$

$\swarrow_{CE}$

where $p_{model}(y^{(t)} | \{x^{(1)}, \ldots, x^{(\tau)}\})$ is given by reading the entry for $y^{(t)}$ from the model's output vector $\hat{y}^{(t)}$

## Loss Function

Total loss is sum of the losses over all the time steps.
So, if $L^{(t)}$ is the negative log likelihood of $y^{(t)}$ given $x^{(1)}, \ldots, x^{(\tau)}$, then

$$
\begin{aligned}
& L(\{x^{(1)}, \ldots, x^{(\tau)}\}, \{y^{(1)}, \ldots, y^{(\tau)}\}) \\
= & \sum_t L^{(t)} \\
= & -\sum_t \log p_{model}(y^{(t)} | \{x^{(1)}, \ldots, x^{(\tau)}\})
\end{aligned}
$$

where $p_{model}(y^{(t)} | \{x^{(1)}, \ldots, x^{(\tau)}\})$ is given by reading the entry for $y^{(t)}$ from the model's output vector $\hat{y}^{(t)}$
Back propagation - right to left - back propagation through time (BPTT)

## A RNN Language Model

$$\hat{y}^{(4)} = P(x^{(5)}|\text{the students opened their})$$

$$h^{(1)} = \tanh\left(Wh^{(0)} + Ux^{(i)} + b\right)$$

**output distribution**

$$\hat{y}^{(t)} = \text{softmax}\left(V h^{(t)} + b_2\right) \in \mathbb{R}^{|V|}$$

**hidden states**

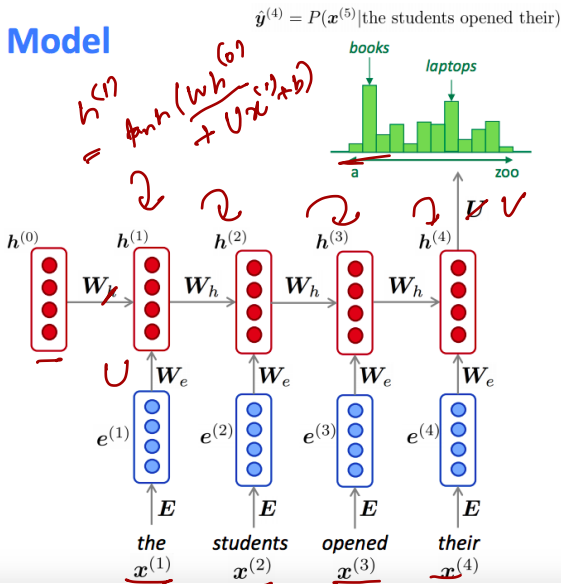$$h^{(t)} = \sigma\left(W_h h^{(t-1)} + W_e e^{(t)} + b_1\right)$$

$h^{(0)}$ is the initial hidden state

**word embeddings**

$$e^{(t)} = E x^{(t)}$$

**words / one-hot vectors**

$$x^{(t)} \in \mathbb{R}^{|V|}$$