1-column per word

1-hot

barley

$n \times d$

$v^{th}$ barley

V

d

1

$v^{th}$ barley

V

$v_1$

d

x

1

V

copied

d

hidden layer

$v \times 1$

common

$(V_O^{apple})^T (V_{P_1}^{banks})$

$d$

$(V_O^{apple})^T$

apple

$d$

$C_1$

$C_3$

$W_2$

money

$V$

$C_4$

$V_I^I$

banks

$V$

$C_2$

$V_1$

$V_P$
banks

$V$

logits

hidden

output

$(V_O^{money})^T (V_P^{banks})$

loss

$- \log P(GT)$

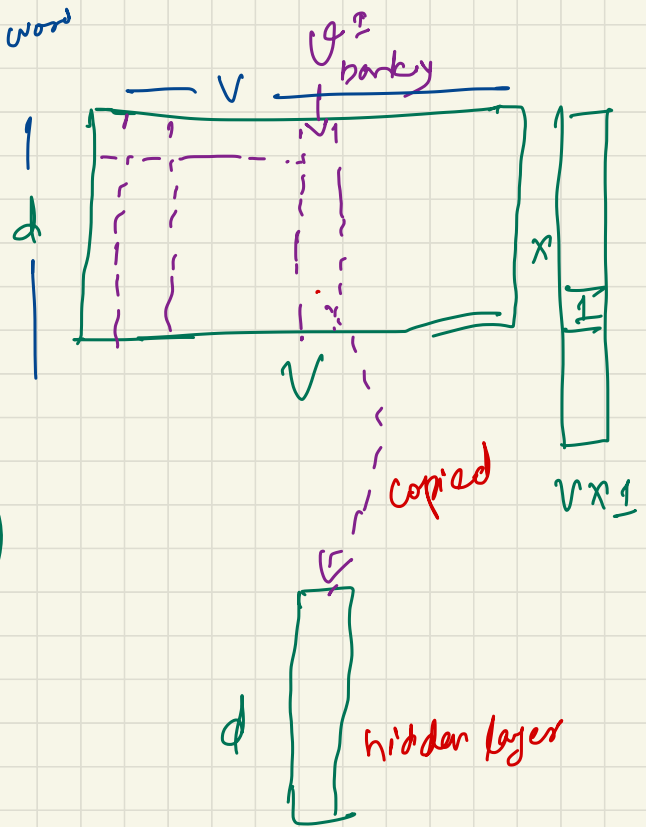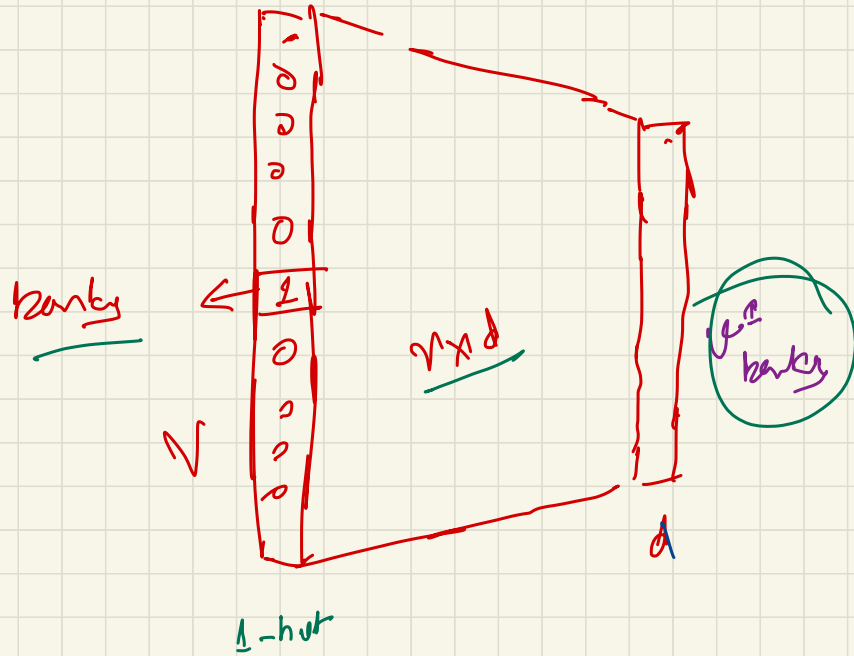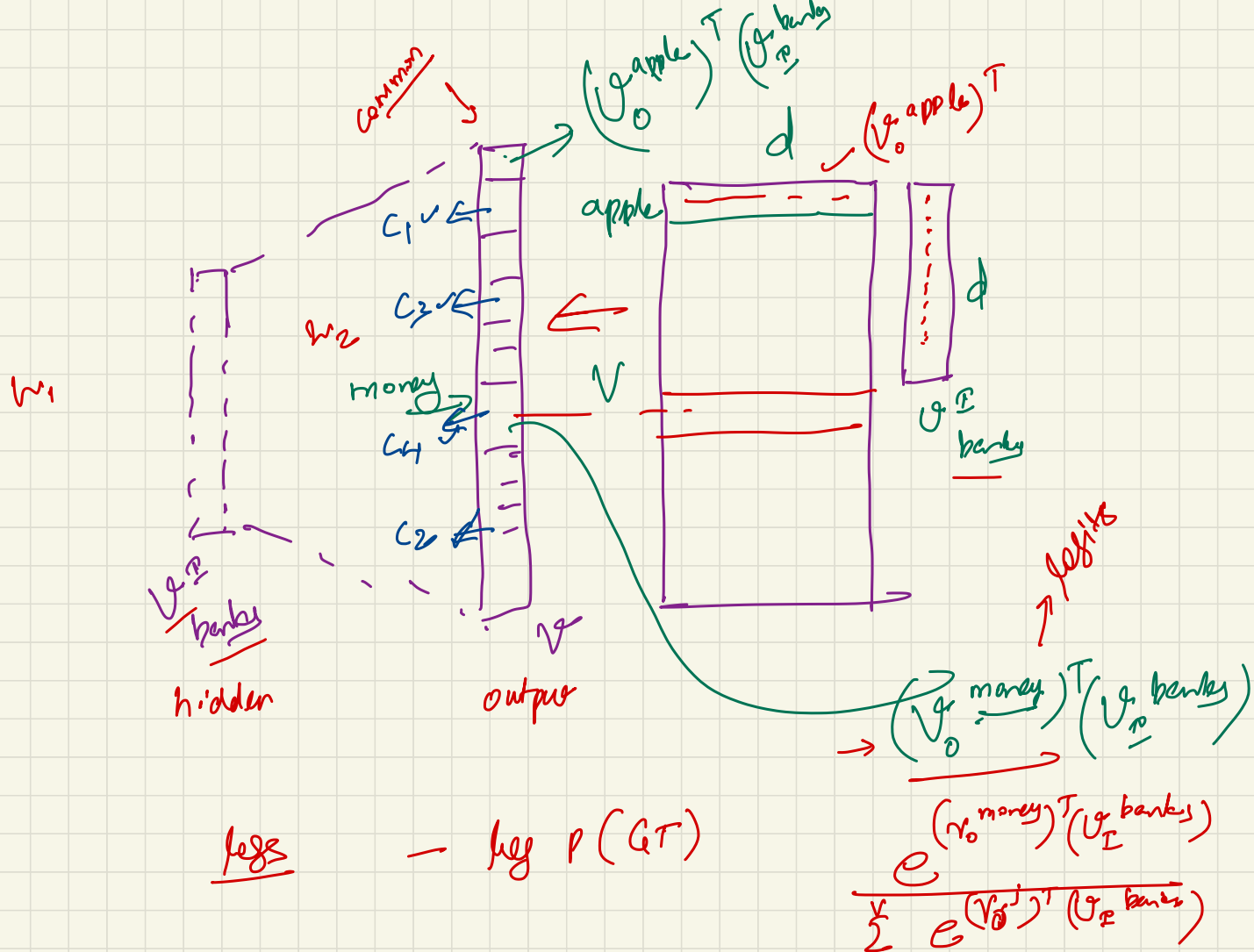$$\frac{e^{(V_O^{money})^T (V_I^{banks})}}{\sum\limits_{V} e^{(V_O^{j})^T (V_I^{banks})}}$$

— Language Models,

— n-gram LMs

issues?

— Parameters depends on window length

$V^k$ (for $k$-prev. words)

# Word2Vec: objective function

We want to minimize the loss function:

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

*center*

*context*

**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 → *Input*
- $u_w$ when $w$ is a context word → *output*

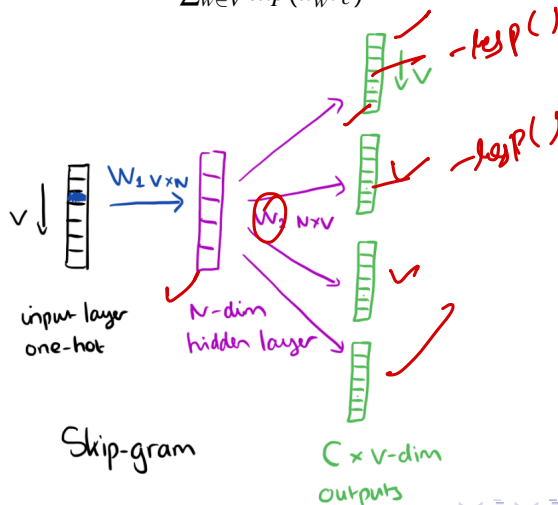Then, for a center word $c$ and a context word $o$

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

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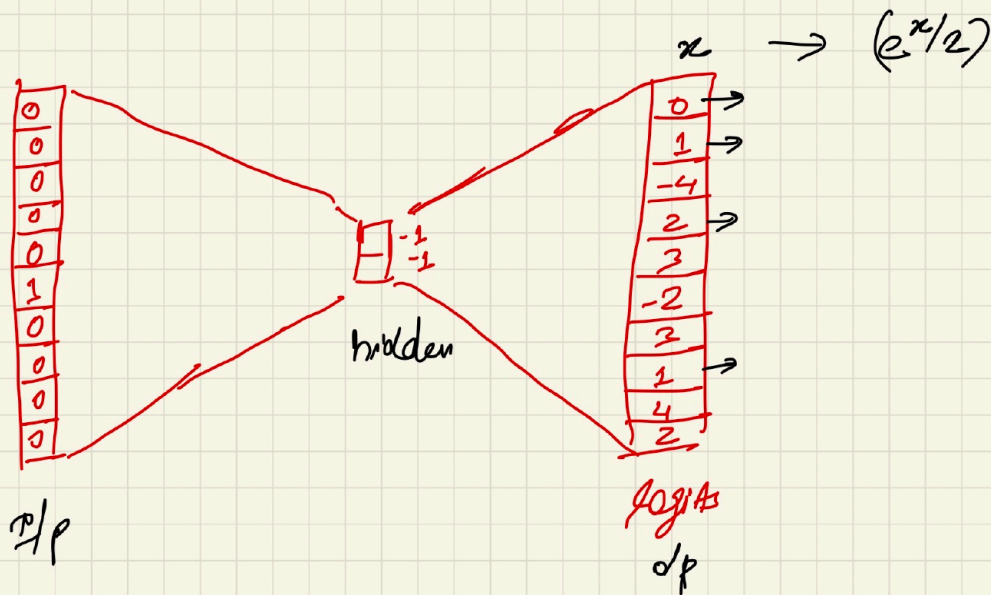


Skip-gram

Suppose you are computing the word vectors using Skip-gram architecture. You have 10 words in your vocabulary, $\{hi, you, they, are, am, how, why, there, who, what\}$ in that order and suppose you have the window, 'hi there how are you' in your corpora. You use this window with 'how' as the center word and two words before and after the center word as your context. Also, suppose that for each word, you have 2-dim in and out vectors, which have the following value at this point given as follows:

*(handwritten annotations: "center", "window", "V = 10")*

| Word | In-vector | Out-vector |
|------|-----------|------------|
| hi | (1, -1) | (-2, 2) |
| you | (-1, 2) | (1, -2) |
| they | (-2, -2) | (2, 2) |
| are | (1, 1) | (-1, -1) |
| am | (2, 1) | (-2, -1) |
| how | (-1, -1) | (1, 1) |
| why | (1, 2) | (-1, -2) |
| there | (2, -1) | (-2, 1) |
| who | (2, 2) | (-2, -2) |
| what | (1, 1) | (-1, -1) |

*(handwritten values next to Out-vector column: 0, 3, -4, 2, 3, -2, 3, 1, 4, 2)*

Table 1: In and Out representations for words

What will the values at the input, hidden and output layer as per the Skip-gram architecture? What would be the total loss for this window?

$x \rightarrow (e^x/2)$

Input column: 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0

hidden: -1, 1

Output column (logits): 0, 1, -4, 2, 3, -2, 3, 1, 4, 2

50-d
25-d
300-d embeddings

I/p

logits

o/p

$$Z = e^0 + 2 \cdot e^1 + 2 \cdot e^2 + 2 \cdot e^3 + e^4 + e^{-2} + e^{-4}$$

$$\text{Total loss} = -\left[ \log\left(\frac{e^0}{z}\right) + \log\left(\frac{e^1}{z}\right) + \log\left(\frac{e^2}{z}\right) + \log\left(\frac{e}{z}\right) \right]$$

### *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.
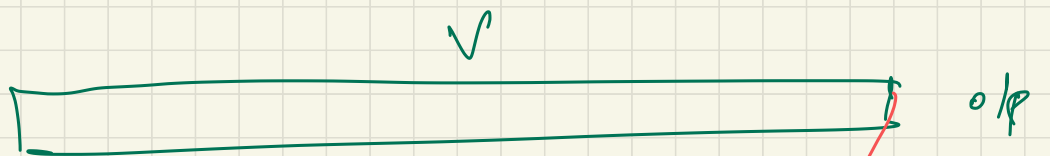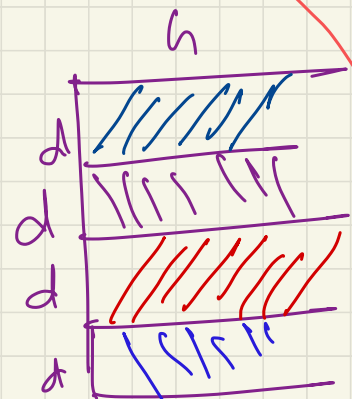
### *Compute the loss*

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 $x$?
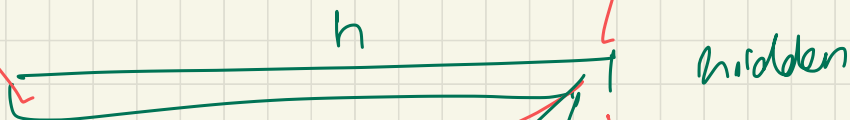
# Homework

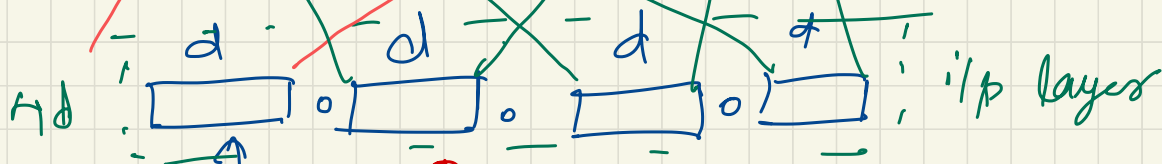- Compute partial derivative of the loss with respect to $v_c$

$$\frac{\partial J}{\partial v_c}$$

$V$

o/p

discrete

$h$

hidden

$d \ll V$

$G$

$d$

$d$

$d$

$d$

$A\phi$

$d$ · · $d$ - - $d$ - - $\phi$

; i/p layer

o    o    o

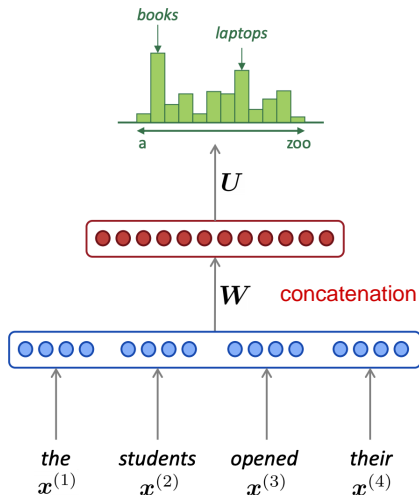the    students    opened    their

students    have    opened    their

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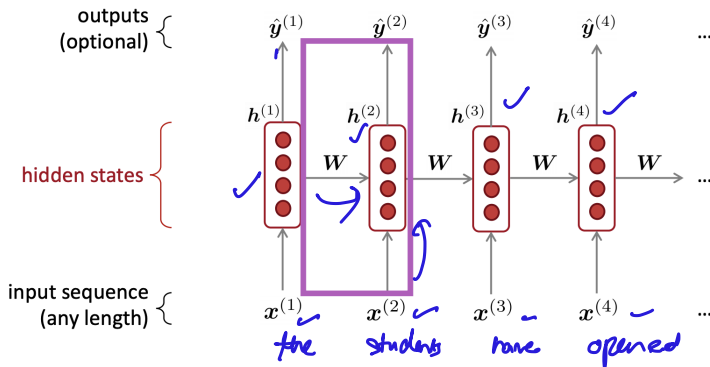- No sparsity problem
- Don't need to store all observed $n$-grams

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*



$U$

$W$    concatenation

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

books

laptops

a    zoo

# Recurrent Neural Networks



outputs (optional): $\hat{\boldsymbol{y}}^{(1)}$, $\hat{\boldsymbol{y}}^{(2)}$, $\hat{\boldsymbol{y}}^{(3)}$, $\hat{\boldsymbol{y}}^{(4)}$ ...

hidden states: $\boldsymbol{h}^{(1)}$, $\boldsymbol{h}^{(2)}$, $\boldsymbol{h}^{(3)}$, $\boldsymbol{h}^{(4)}$

$\boldsymbol{W}$, $\boldsymbol{W}$, $\boldsymbol{W}$, $\boldsymbol{W}$ ...

input sequence (any length): $\boldsymbol{x}^{(1)}$, $\boldsymbol{x}^{(2)}$, $\boldsymbol{x}^{(3)}$, $\boldsymbol{x}^{(4)}$ ...

the    students    have    opened

## Core Idea

Apply the same weights repeatedly!
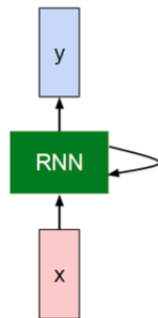
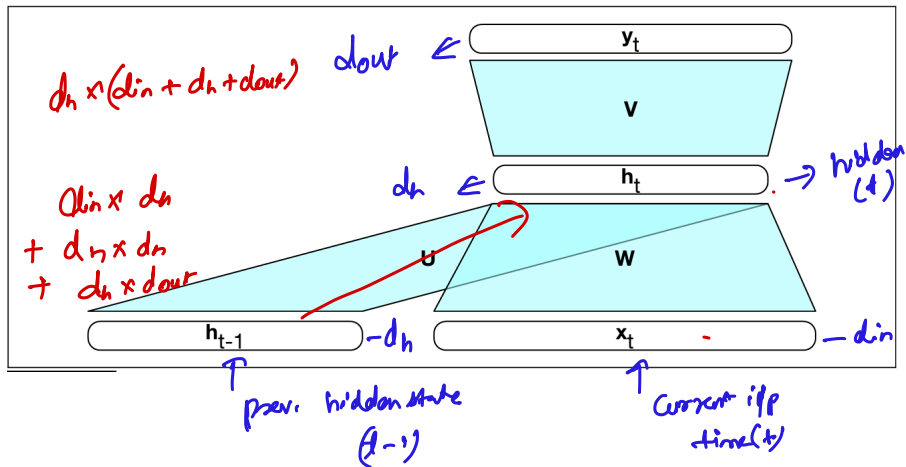We can process a sequence of vectors $x$ by applying a recurrence formula at each step:

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

new state

some function with parameters W

old state

input vector at some time step

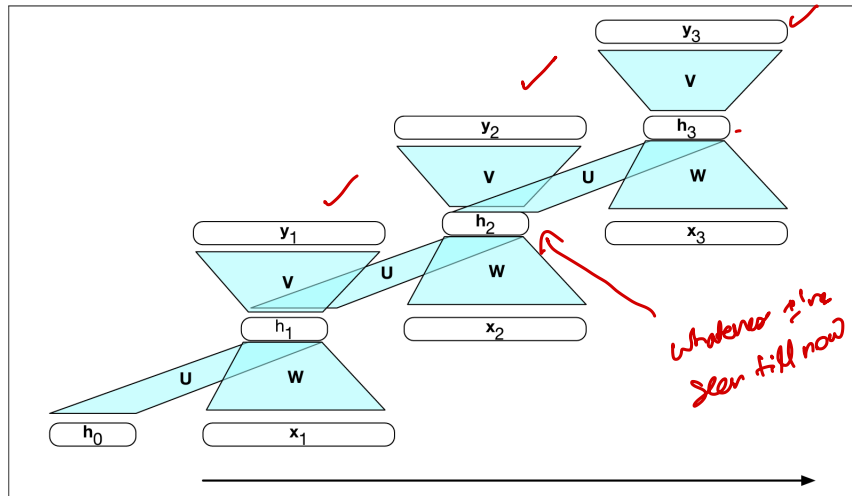Notice: the same function and the same set of parameters are used at every time step.

$d_h \times (d_{in} + d_h + d_{out})$

$d_{out}$

$y_t$

V

$d_h$

$h_t$ → hidden (t)

$d_{in} \times d_h$
$+ d_h \times d_h$
$+ d_h \times d_{out}$

U

W

$h_{t-1}$ — $d_h$

prev. hidden state (t-1)

$x_t$ — $d_{in}$

current i/p time (t)

## Forward Propagation

$$h_t = g(Uh_{t-1} + Wx_t)$$
$$y_t = softmax(Vh_t)$$

- Let the dimensions of the input, hidden and output be $d_{in}$, $d_h$ and $d_{out}$, respectively
- The three parameter matrices: $W : d_h \times d_{in}$, $U : d_h \times d_h$, $V : d_{out} \times d_h$

highest—prob (students)

loop

opened

loss

$h_T$

$y_f$

boy students

$h_T$

$V$

$V$

$h_0$ $U$ $\rightarrow$ $h_1$ $U$ $\rightarrow$ $h_2$

$U, W, V$

$W$

$W$ $\checkmark$

BPTT

back-propagation through time

$\checkmark$

$\rightarrow$ Teacher forcing

(only during training)

$x_1$

the    students    opened    their

- To train RNN LM, we use self-supervision (or self-training)
- We take a corpus of text as training material
- At each time step $t$, we ask the model to predict the next word

*Why is it called self-supervision?*

- We do not add any gold data, the natural sequence of words is its own supervision!
- We simply train the model to minimize the error in predicting the true next word in the training sequence

# Training an RNN language model

argmax

$z_1$

$v\_dim$

Season

greedy
decoding

Prob.
dist

random
sampling

$h_0$

$h_1$

$x_1$

< Start >

$x_2$

The