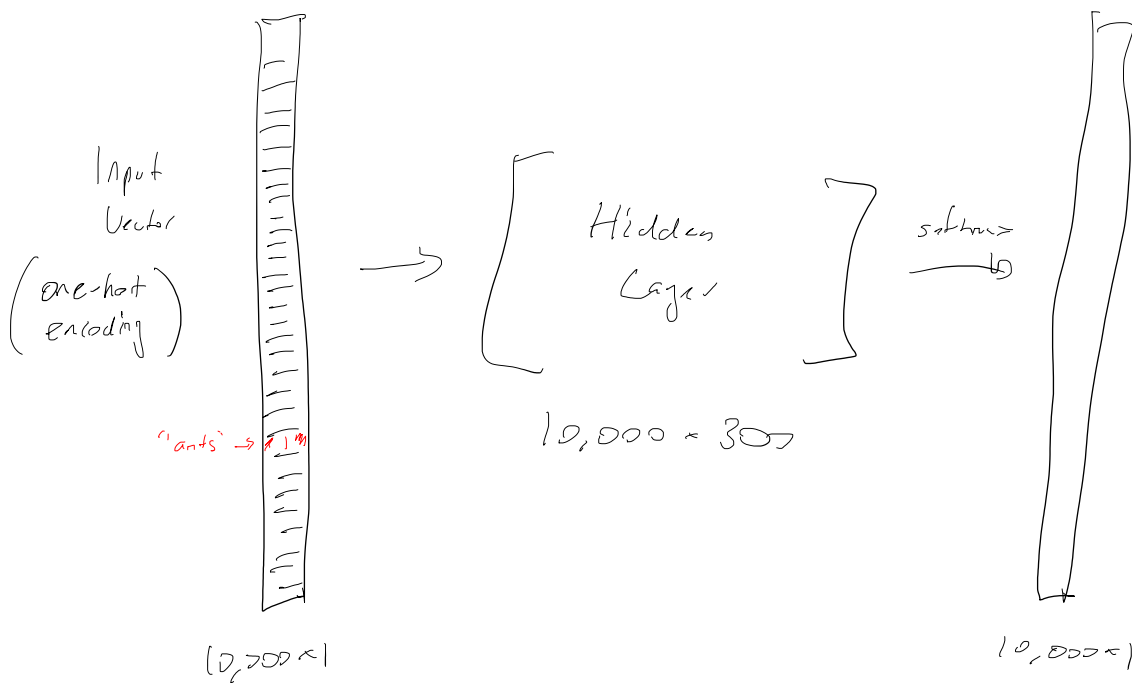


Word 2 Vec

Skip-gram

- Similar task as using encoder/decoder for unsupervised learning + strip away output layer, use just hidden layer
- Give model a "fake task", i.e. given a word, estimate the probability of being a context word, for each word in the vocabulary

↳ i.e. word is within a window



$$\begin{bmatrix} 0 & 0 & 0 & 1 & 0 \end{bmatrix} \times \begin{bmatrix} 17 & 24 & 1 \\ 23 & 5 & 7 \\ 4 & 6 & 13 \\ 10 & 12 & 19 \\ 11 & 18 & 25 \end{bmatrix} = \begin{bmatrix} 10 & 12 & 19 \end{bmatrix}$$

↙ Softmax Classifier

$$\text{Softmax} : \frac{e^{\theta x}}{\sum e^{\theta x}}$$

So: Input is a $1 \times 10,000$ one-hot vector.
 Hidden layer is $10,000 \times 300$ hidden layer.

The product is a 1×300 vector.

This is input into a softmax classifier

(So each output neuron has a weight vector, θ_i)
 then calculates $\exp(\theta_i X)$
 \uparrow output, 1×300 vector

then to normalize, $\frac{\exp(\theta_i X)}{\sum \exp(\theta_i X)}$

CBOW (Continuous bag of words)

• Predict center word from "context"

the cat jumped over the fence.

• Let

• C = window size

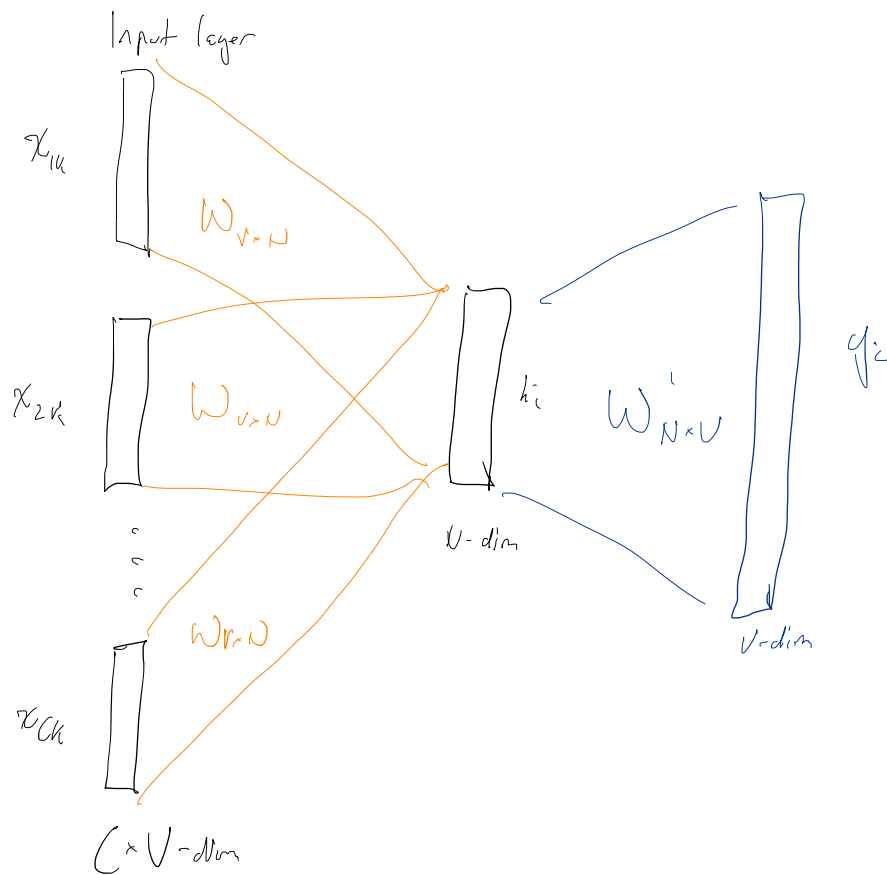
• V = vocabulary size

• h : hidden layer, n -dimensional vector (n is arbitrary)

• W : $V \times n$ weight matrix

• W' : $n \times V$ weight matrix

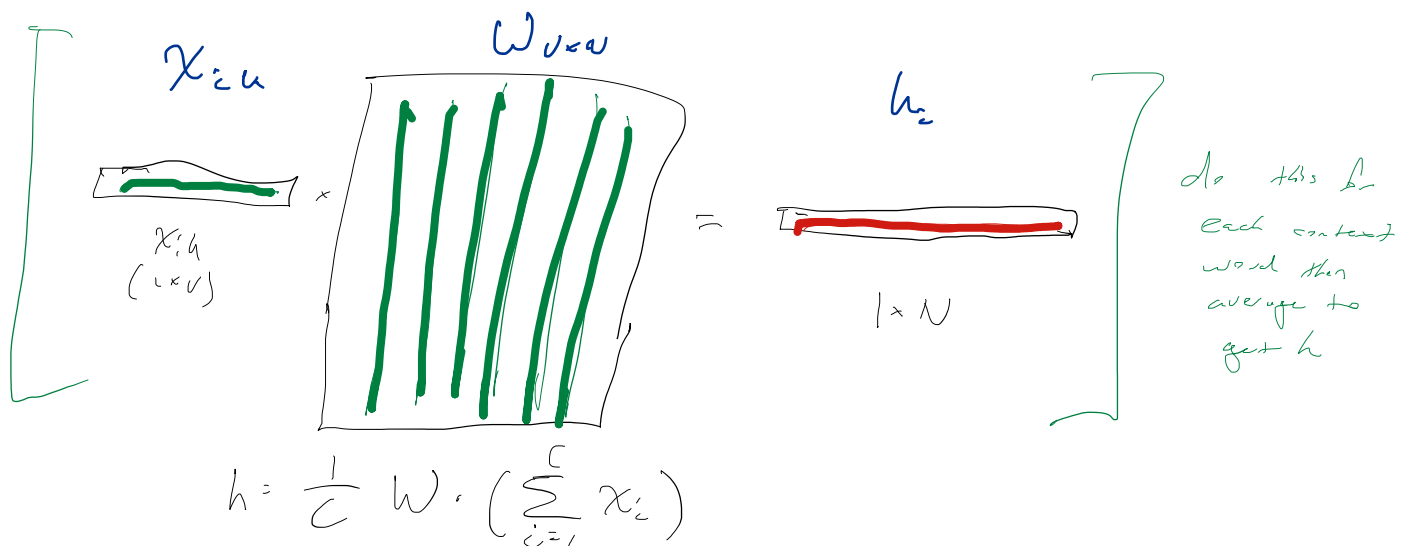
• (skip-gram model, reversed)



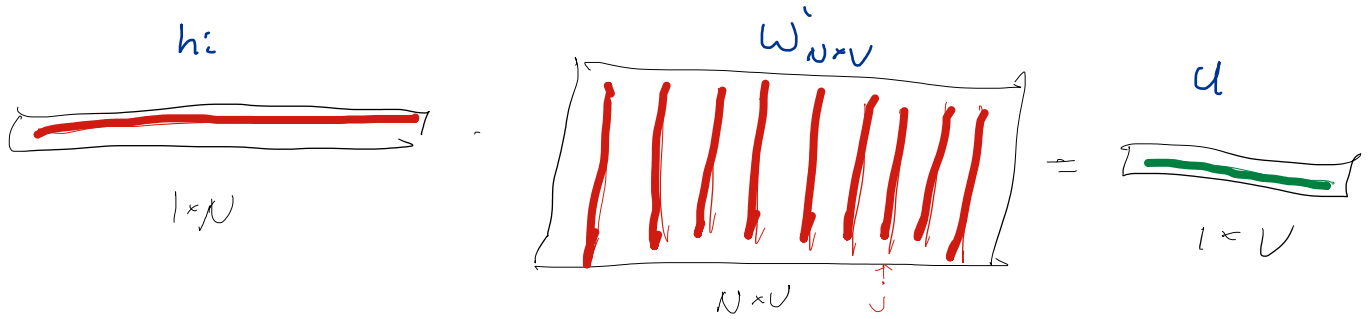
① $\{x_{1u}, x_{2u}, \dots, x_{cu}\}$ are one-hot encoded vectors

② Concatenated to form $C \times V$ matrix

③ Multiply with weights $W_{V \times N}$ to get $C \times N$ matrix, which is averaged to get a $1 \times N$ vector.



(4) Multiply the hidden layer vector ($1 \times N$) with W' ($N \times V$)



Remember:

$$(B_0 \ B_1) \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix}$$

$$= (a_{11} B_0 + a_{21} B_1, a_{12} B_0 + a_{22} B_1)$$

So $u_j = v_{w_j}' \cdot h$

↑ hidden layer

↑ j -th column from matrix $W'_{N \times V}$

↑ pass u_j through softmax

$$y_j = p(w_{qj} | w_1, \dots, w_c) = \frac{\exp(u_j)}{\sum_{j=1}^V \exp(u_j)}$$

softmax

Softmax:
e.g.

$$y = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} = \begin{bmatrix} 2.0 \\ 1.0 \\ 0.1 \end{bmatrix} \quad \frac{e^{y_1}}{\sum_j e^{y_j}} = \frac{e^2}{e^2 + e^1 + e^{0.1}} \approx 0.7$$

