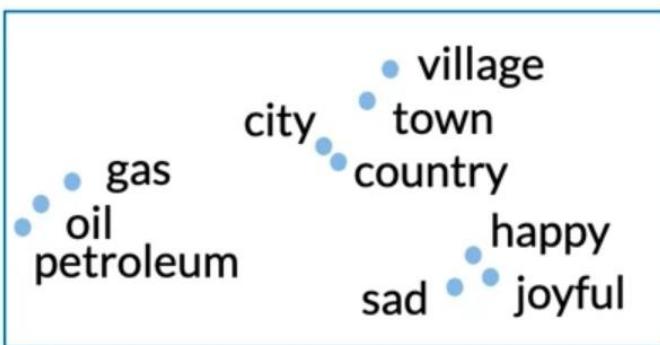


Some basic applications of word embeddings



Semantic analogies
and similarity



Sentiment analysis

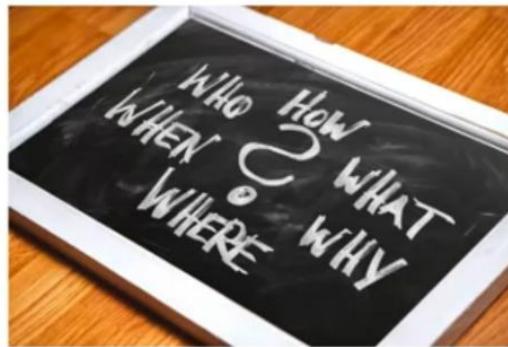


Classification of
customer feedback

Advanced applications of word embeddings



Machine translation



Information extraction

Learning objectives

Prerequisite: neural networks

- Identify the key concepts of word representations
- Generate word embeddings
- Prepare text for machine learning
- Implement the continuous bag-of-words model

Integers

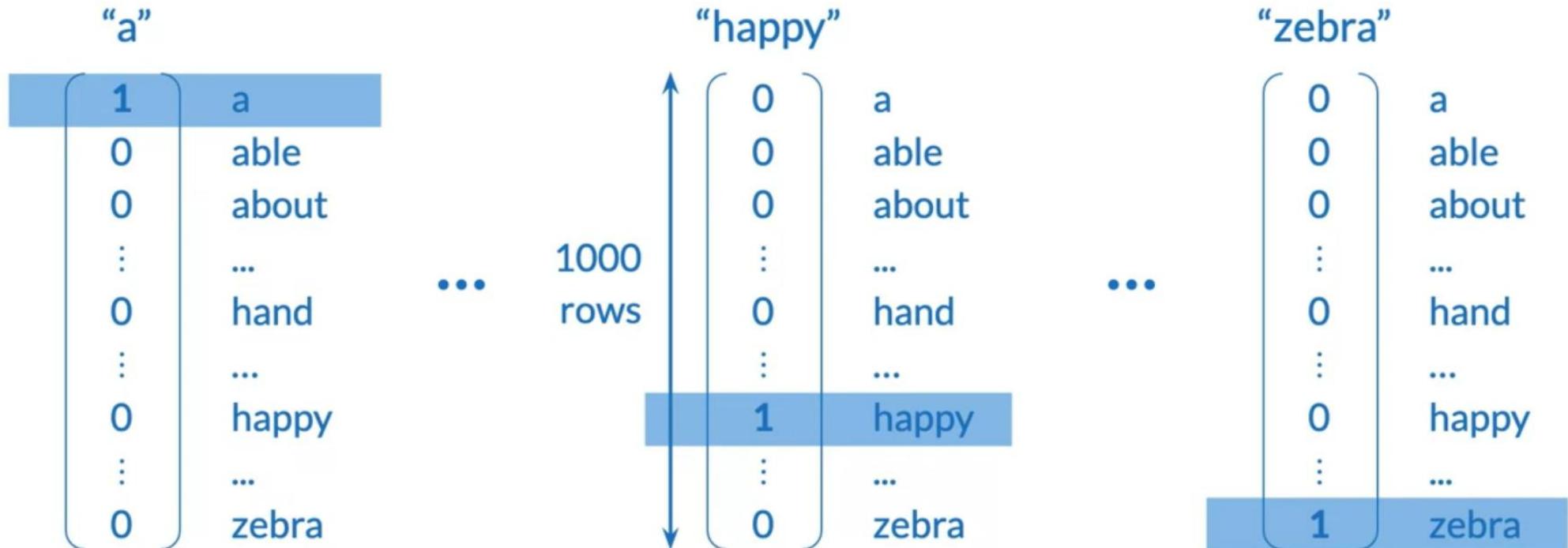
Word	Number
a	1
able	2
about	3
...	...
hand	615
...	...
happy	621
...	...
zebra	1000

Integers

- + Simple
- Ordering: little semantic sense

hand < happy < zebra
615 621 1000
?! ?!

One-hot vectors

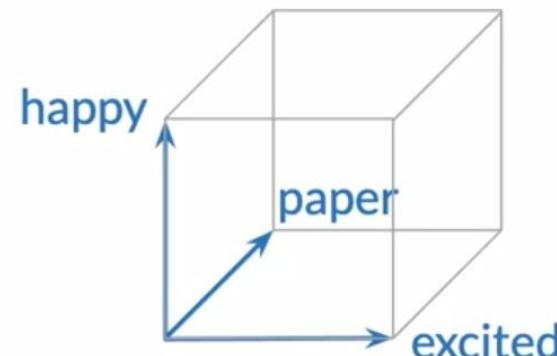
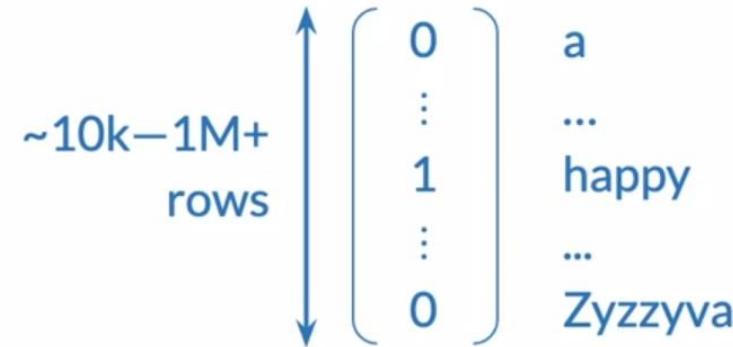


One-hot vectors

Word	Number	"happy"		
a	1	1	0	a
able	2	2	0	able
about	3	3	0	about
...
hand	615	615	0	hand
...
happy	621	621	1	happy
...
zebra	1000	1000	0	zebra

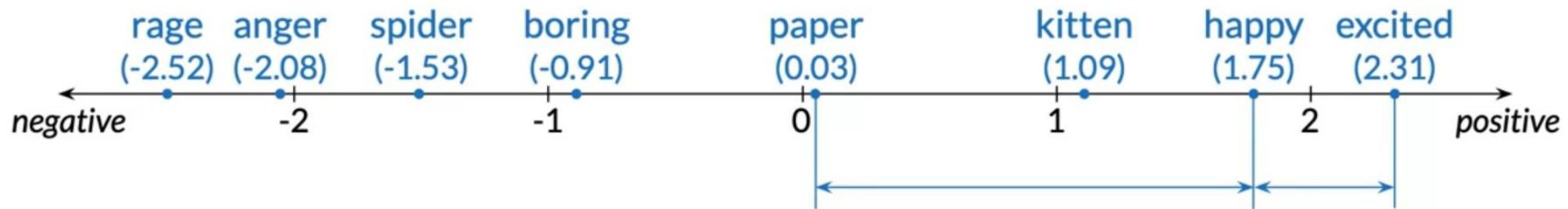
One-hot vectors

- + Simple
- + No implied ordering
- Huge vectors
- No embedded meaning



$d(\text{paper}, \text{excited})$
 $= d(\text{paper}, \text{happy})$
 $= d(\text{excited}, \text{happy})$

Meaning as vectors

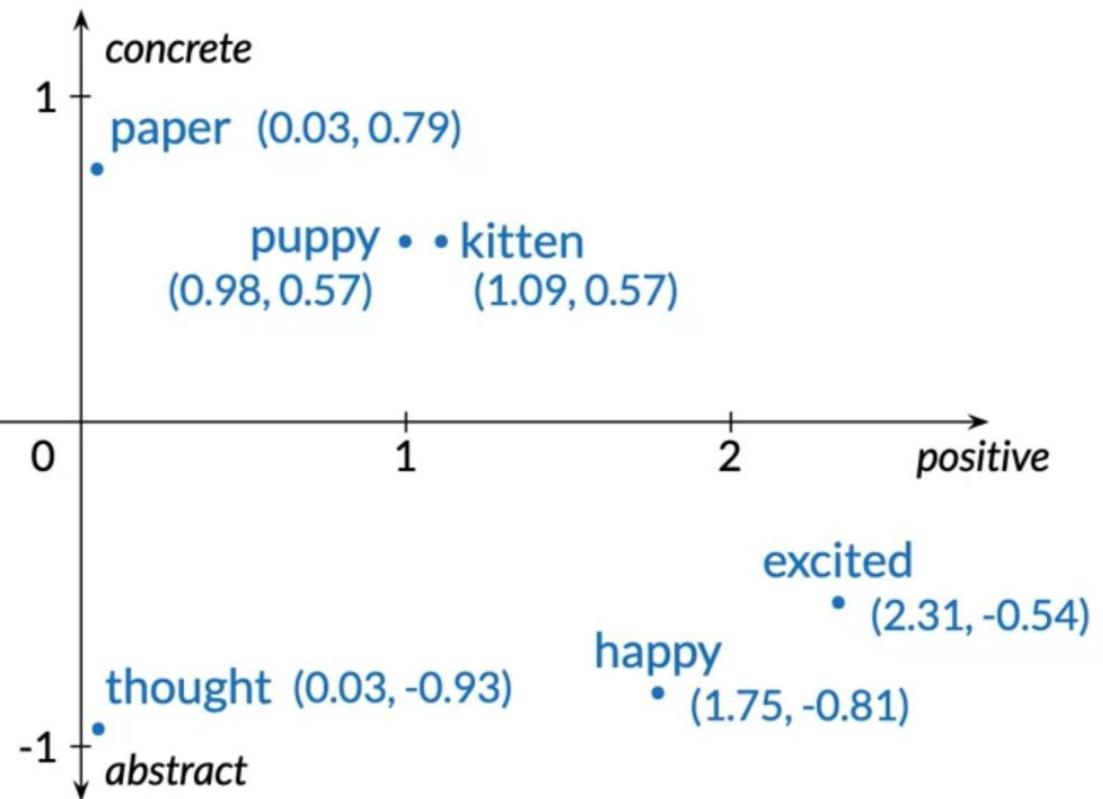


Meaning as vectors

negative -2 -1 0 positive

spider (-1.53, 0.41)
snake (-1.53, 0.41)

rage (-2.52, -0.54)
anger (-2.08, -0.71)
boring



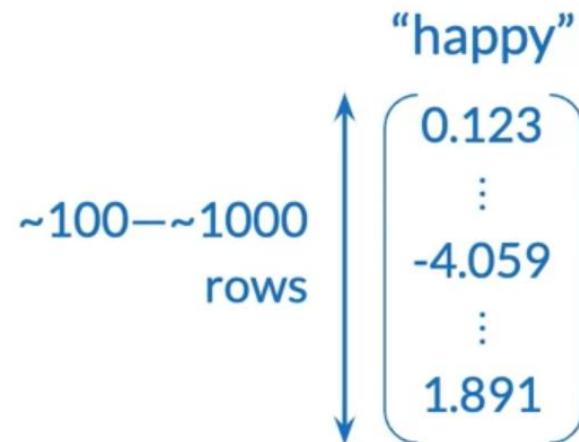
Word embedding vectors

- + Low dimension
- + Embed meaning
 - o e.g. semantic distance

forest \approx tree forest \neq ticket

- o e.g. analogies

Paris:France :: Rome:?



Terminology

integers

word vectors

one-hot vectors

word embedding vectors

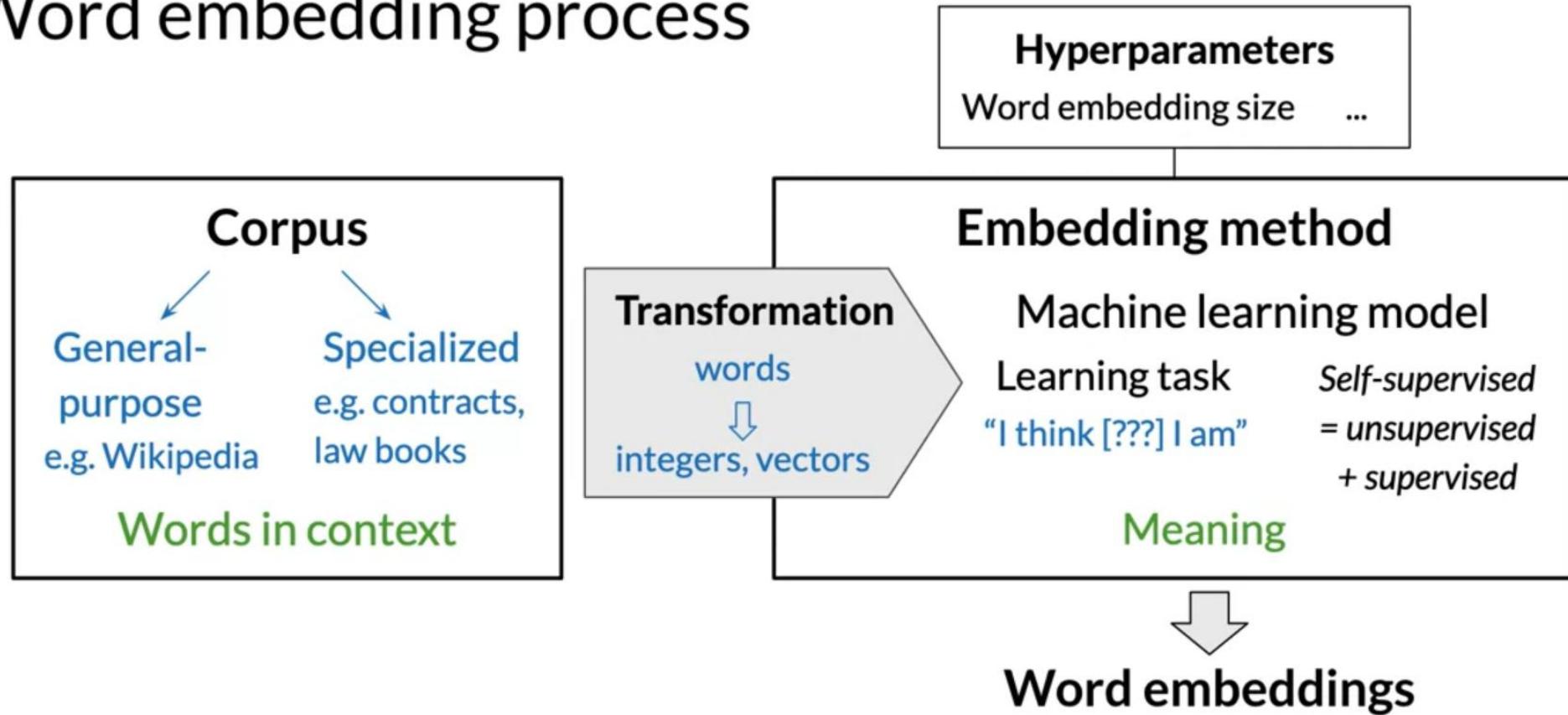
“word vectors”

word embeddings

Summary

- Words as integers
- Words as vectors
 - One-hot vectors
 - Word embedding vectors
- Benefits of word embeddings for NLP

Word embedding process



Basic word embedding methods

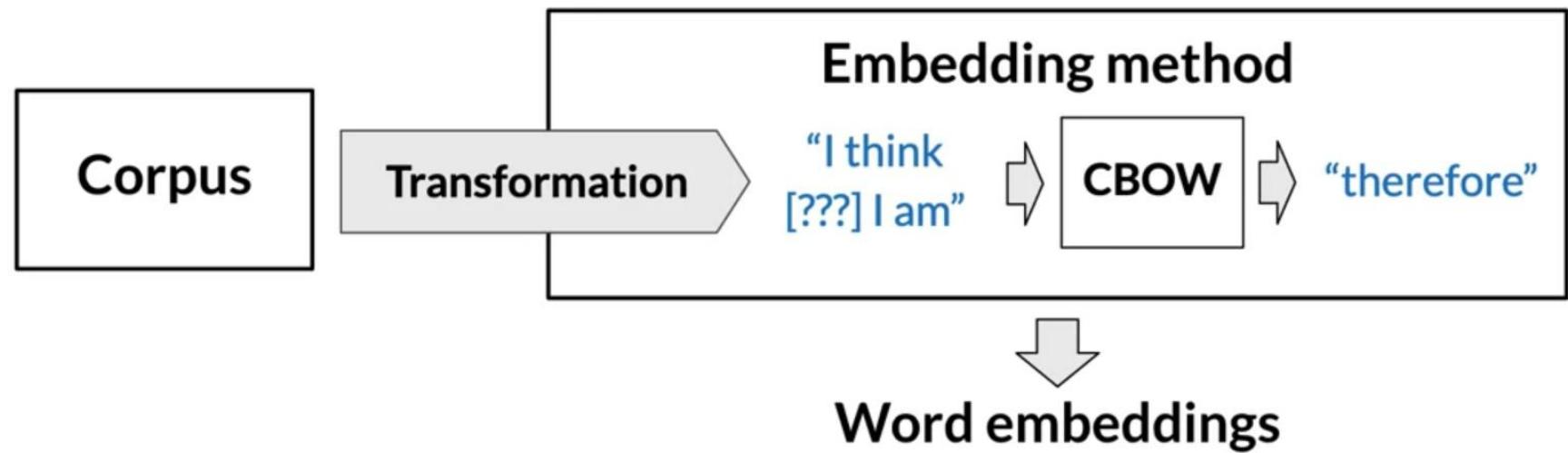
- word2vec (Google, 2013)
 - Continuous bag-of-words (CBOW)
 - Continuous skip-gram / Skip-gram with negative sampling (SGNS)
- Global Vectors (GloVe) (Stanford, 2014)
- fastText (Facebook, 2016)
 - Supports out-of-vocabulary (OOV) words

Advanced word embedding methods

Deep learning, contextual embeddings

- BERT (Google, 2018)
 - ELMo (Allen Institute for AI, 2018)
 - GPT-2 (OpenAI, 2018)
- 
- Tunable pre-trained
models available

Continuous bag-of-words word embedding process





Center word prediction: rationale

The little _____ ? is barking



dog

puppy

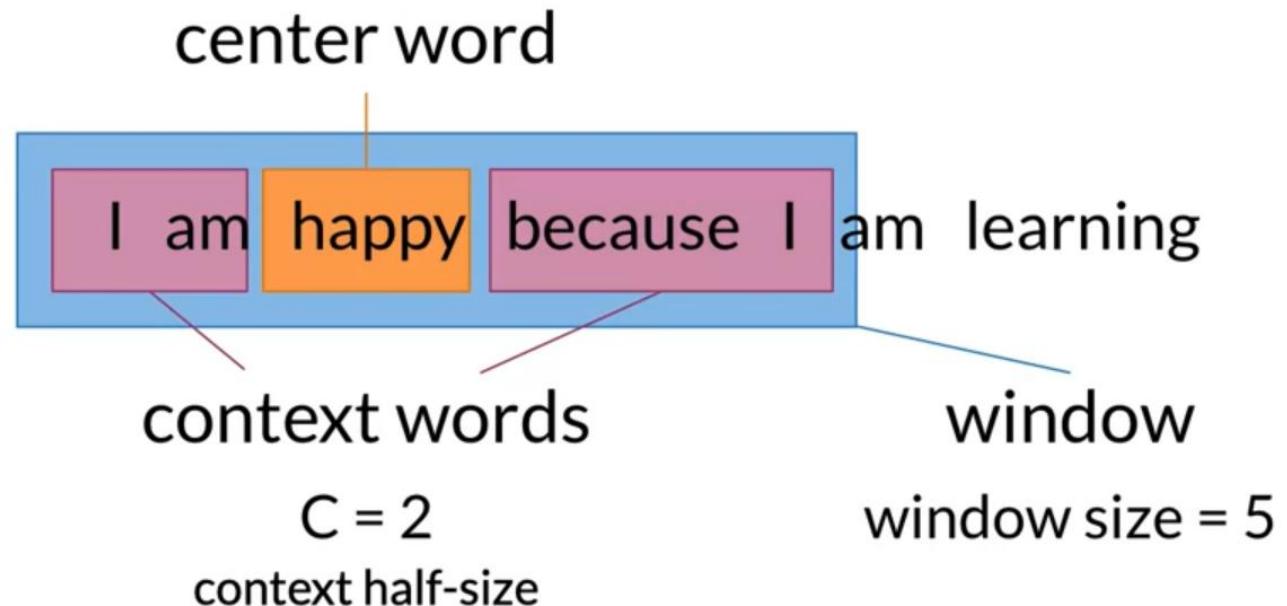
hound

terrier

...

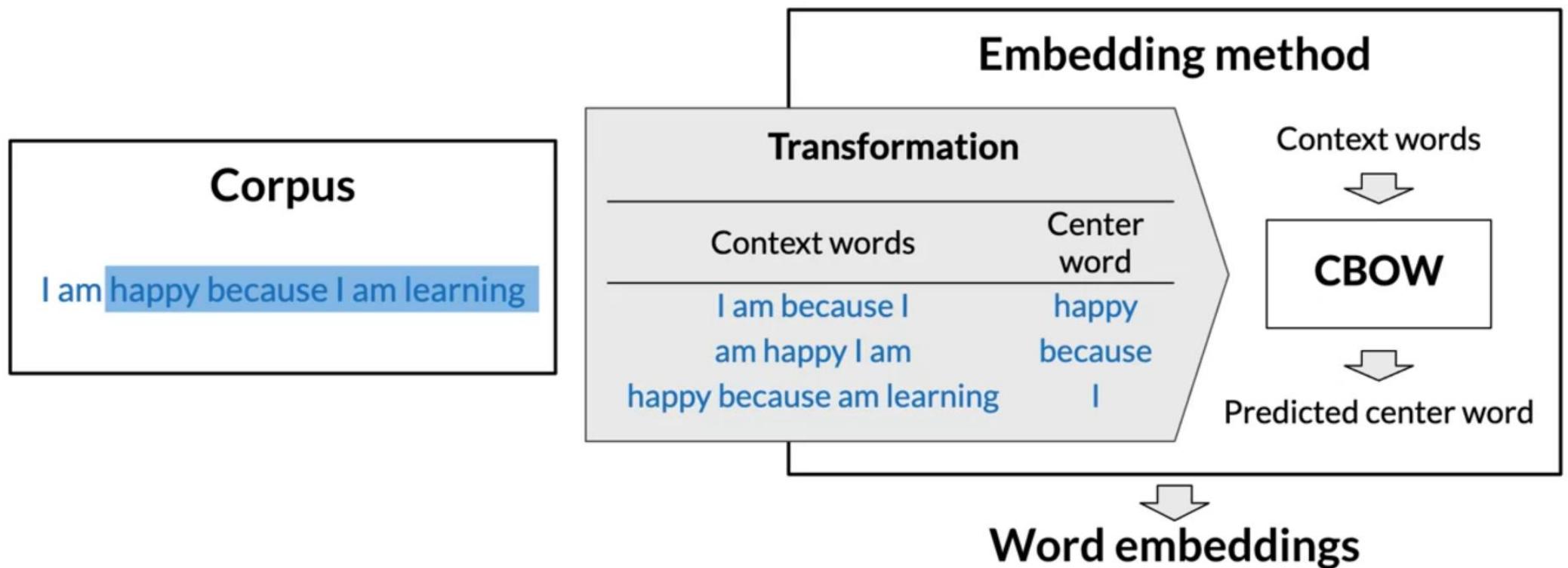


Creating a training example

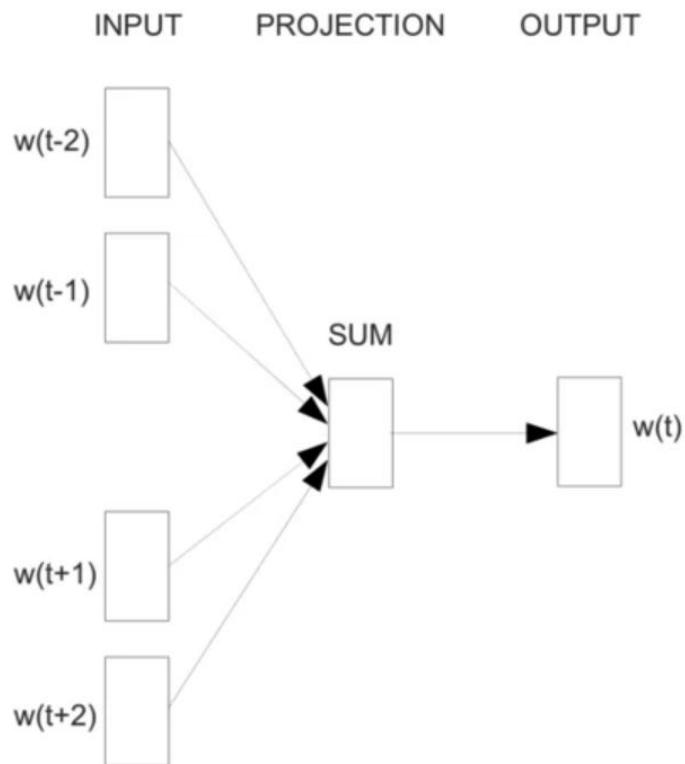




From corpus to training



CBOW in a nutshell



Source: Mikolov, T., Chen, K., Corrado, G.S., & Dean, J. (2013).
[Efficient Estimation of Word Representations in Vector Space](#)

Cleaning and tokenization matters

- Letter case “The” == “the” == “THE” → lowercase / uppercase
- Punctuation , ! . ? → . “ ‘ « » ’ ” → Ø ... !! ??? → .
- Numbers 1 2 3 5 8 → Ø 3.14159 90210 → as is / <NUMBER>
- Special characters Ⓛ \$ € § ¶ ** → Ø
- Special words 😊 #nlp → :happy: #nlp

Example in Python: corpus

Who ❤️ "word embeddings" in 2020? I do!!!

emoji

punctuation

number

Example in Python: libraries

```
# pip install nltk
# pip install emoji

import nltk
from nltk.tokenize import word_tokenize
import emoji

nltk.download('punkt') # download pre-trained Punkt tokenizer for English
```

Example in Python: code

```
corpus = 'Who ❤️ "word embeddings" in 2020? I do!!!'  
  
data = re.sub(r'[,!?;-]+', '.', corpus)
```

→ Who ❤️ "word embeddings" in 2020. I do.

Example in Python: code

```
corpus = 'Who ❤️ "word embeddings" in 2020? I do!!!'

data = re.sub(r'[,!?;-]+', '.', corpus)
data = nltk.word_tokenize(data) # tokenize string to words

→ ['Who', '❤️', '', 'word', 'embeddings', "", 'in', '2020', '.', 'I',
'do', '.']
```

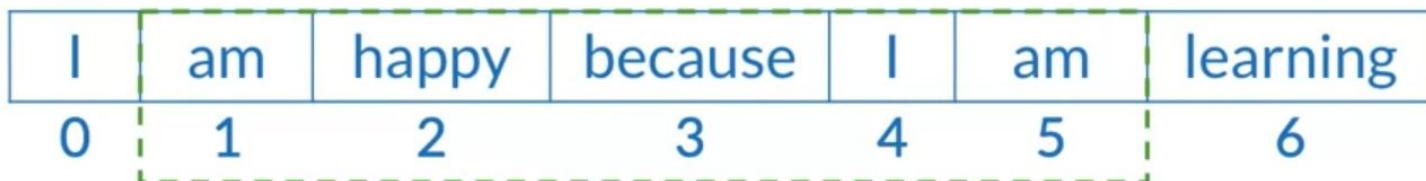
Example in Python: code

```
corpus = 'Who ❤️ "word embeddings" in 2020? I do!!!'

data = re.sub(r'[,!?;-]+', '.', corpus)
data = nltk.word_tokenize(data) # tokenize string to words
data = [ ch.lower() for ch in data
         if ch.isalpha()
         or ch == '.'
         or emoji.get_emoji_regexp().search(ch)
     ]
→ ['who', '❤️', 'word', 'embeddings', 'in', '.', 'i', 'do', '.']
```

Sliding window of words in Python

```
def get_windows(words, C):
    i = C
    while i < len(words) - C:
        center_word = words[i]
        context_words = words[(i - C):i] + words[(i+1):(i+C+1)]
        yield context_words, center_word
        i += 1
```



Sliding window of words in Python

```
def get_windows(words, C):
    ...
        yield context_words, center_word
```

```
for x, y in get_windows(
    ['i', 'am', 'happy', 'because', 'i', 'am', 'learning'],
    2
):
    print(f'{x}\t{y}')
```

Sliding window of words in Python

```
for x, y in get_windows(  
    ['i', 'am', 'happy', 'because', 'i', 'am', 'learning'],  
    2  
):  
    print(f'{x}\t{y}')
```

```
→ ['I', 'am', 'because', 'I']      happy  
['am', 'happy', 'I', 'am']      because  
['happy', 'because', 'am', 'learning']  I
```

Transforming center words into vectors

Corpus I am happy because I am learning

Vocabulary am, because, happy, I, learning

One-hot vector	am	because	happy	I	learning
am	$\begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$	$\begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{pmatrix}$	$\begin{pmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{pmatrix}$	$\begin{pmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{pmatrix}$	$\begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{pmatrix}$
because	$\begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$	$\begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{pmatrix}$	$\begin{pmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{pmatrix}$	$\begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{pmatrix}$	$\begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{pmatrix}$
happy	$\begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{pmatrix}$	$\begin{pmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{pmatrix}$	$\begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$	$\begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{pmatrix}$	$\begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{pmatrix}$
I	$\begin{pmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{pmatrix}$	$\begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{pmatrix}$	$\begin{pmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{pmatrix}$	$\begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$	$\begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{pmatrix}$
learning	$\begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{pmatrix}$	$\begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{pmatrix}$	$\begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{pmatrix}$	$\begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$	$\begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$

Transforming context words into vectors

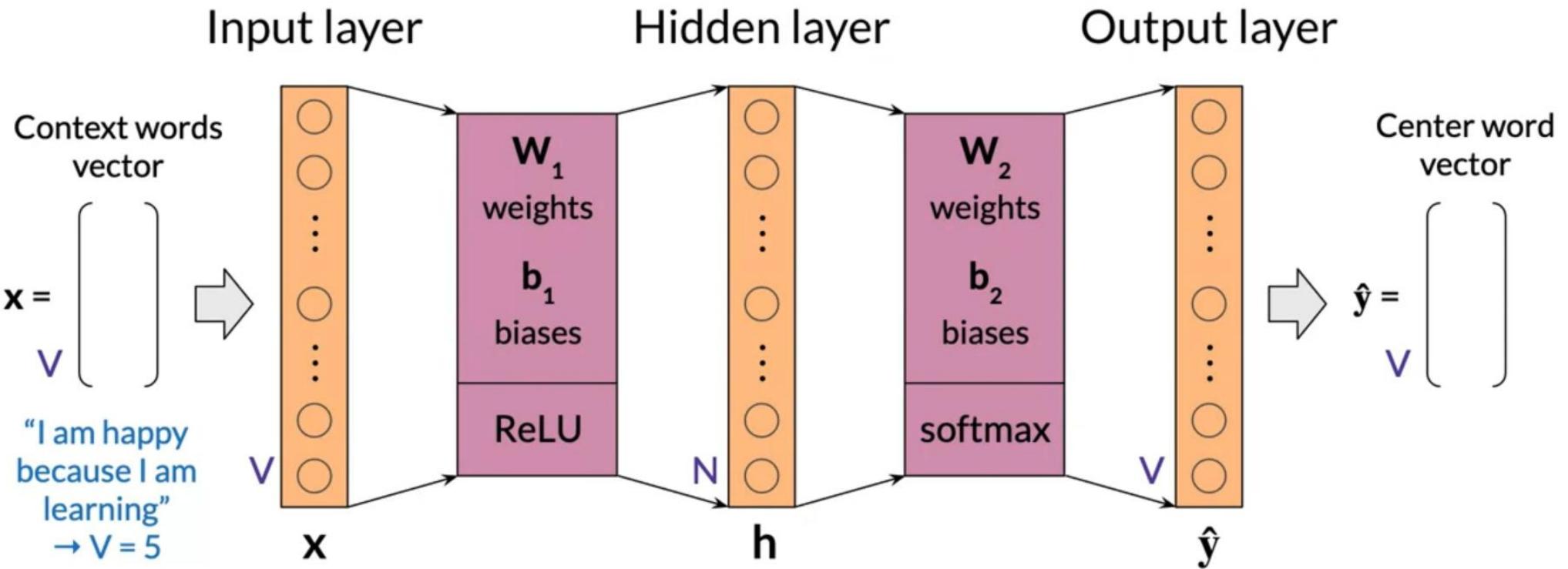
Average of individual one-hot vectors

$$\left[\begin{array}{c} I \\ am \\ because \\ happy \\ I \\ learning \end{array} \right] = \frac{1}{4} \left(\left[\begin{array}{c} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{array} \right] + \left[\begin{array}{c} 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{array} \right] + \left[\begin{array}{c} 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{array} \right] + \left[\begin{array}{c} 0 \\ 0 \\ 0 \\ 1 \\ 1 \\ 0 \end{array} \right] \right)$$

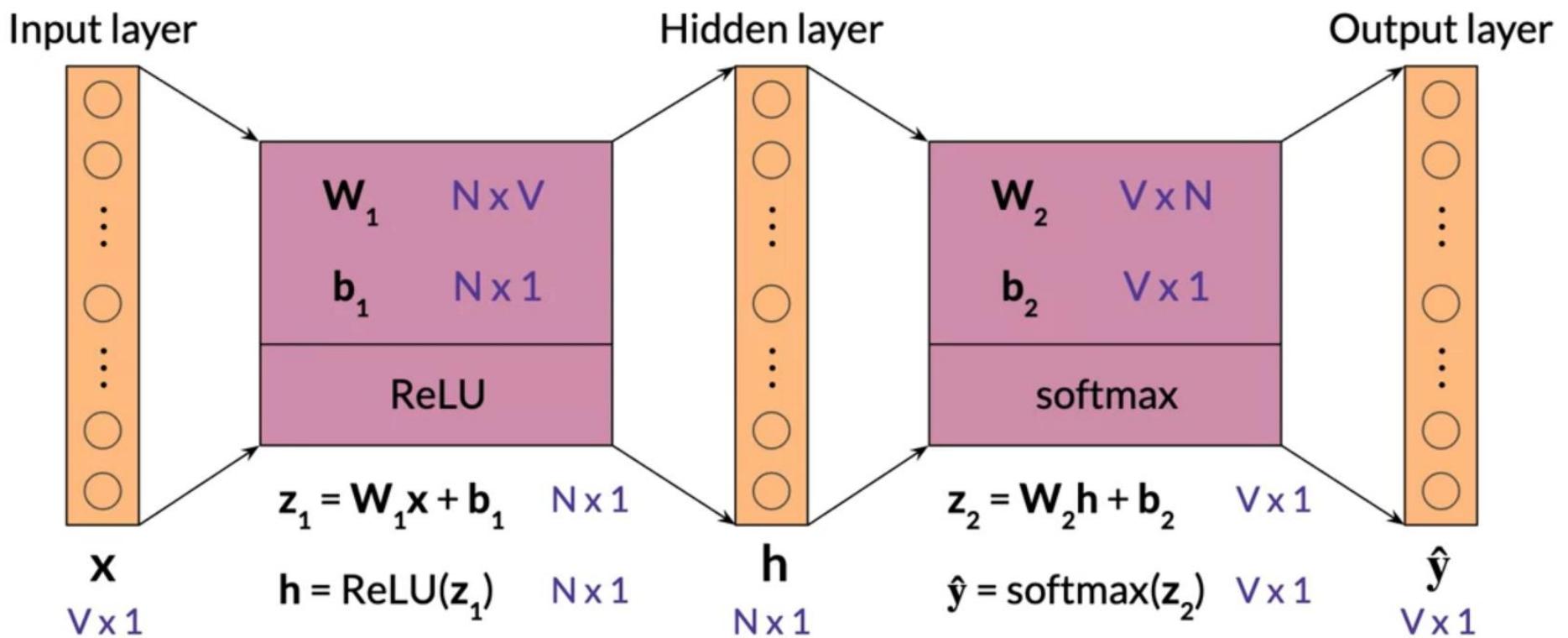
Final prepared training set

Context words	Context words vector	Center word	Center word vector
I am because I	[0.25; 0.25; 0; 0.5; 0]	happy	[0; 0; 1; 0; 0]
am happy I am	[0.5; 0; 0.25; 0.25; 0]	because	[0; 1; 0; 0; 0]
happy because am learning	[0.25; 0.25; 0.25; 0; 0.25]	I	[0; 0; 0; 1; 0]

Architecture of the CBOW model



Dimensions (single input)



Dimensions (single input)

Column vectors

$$z_1 = W_1 x + b_1$$

$$z_1 = \begin{pmatrix} \end{pmatrix} \quad W_1 = \begin{pmatrix} & N \times V \end{pmatrix} \quad x = \begin{pmatrix} \end{pmatrix} \quad b_1 = \begin{pmatrix} \end{pmatrix}$$

$N \times 1$ $V \times 1$ $N \times 1$

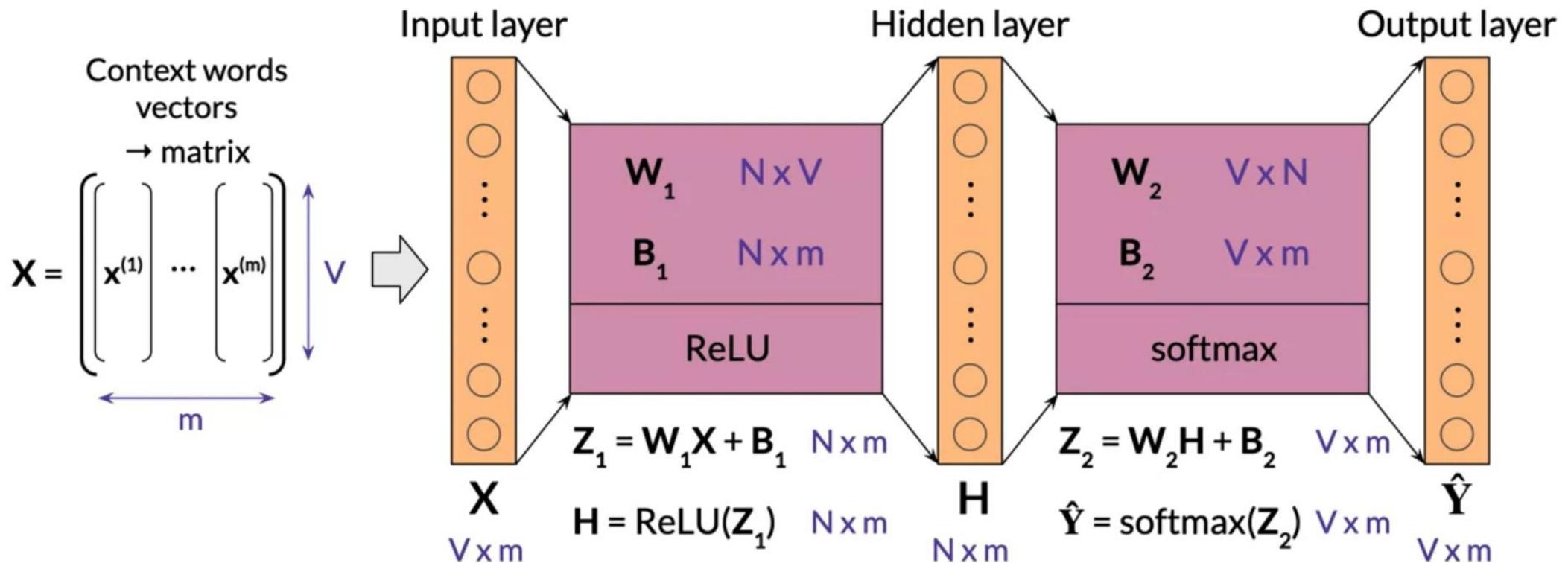
Row vectors

$$z_1 = x W_1^T + b_1$$

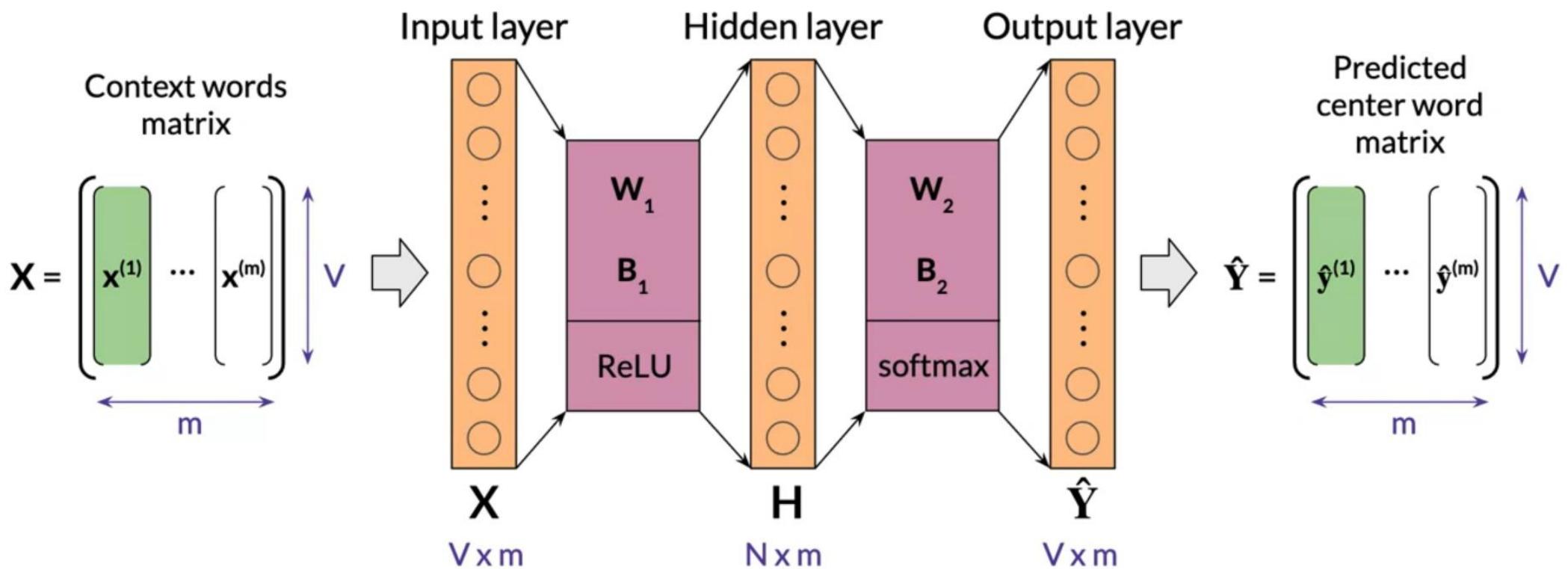
$$b_1 = \begin{pmatrix} 1 \times N \end{pmatrix} \quad W_1 = \begin{pmatrix} & N \times V \end{pmatrix} \quad b_1 = \begin{pmatrix} 1 \times N \end{pmatrix}$$
$$x = \begin{pmatrix} 1 \times V \end{pmatrix}$$

Dimensions (batch input)

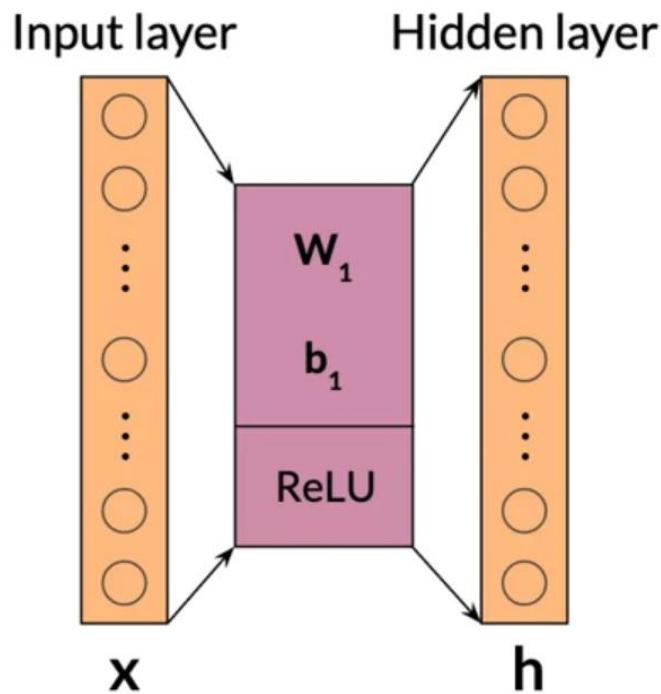
$$\begin{bmatrix} \mathbf{b}_1 \end{bmatrix} \rightarrow \mathbf{B}_1 = \begin{bmatrix} \mathbf{b}_1 \\ \vdots \\ \mathbf{b}_1 \end{bmatrix}_{m \times N} \quad \text{broadcasting}$$



Dimensions (batch input)



Rectified Linear Unit (ReLU)



$$z_1 = W_1 x + b_1$$

$$h = \text{ReLU}(z_1)$$

z_1

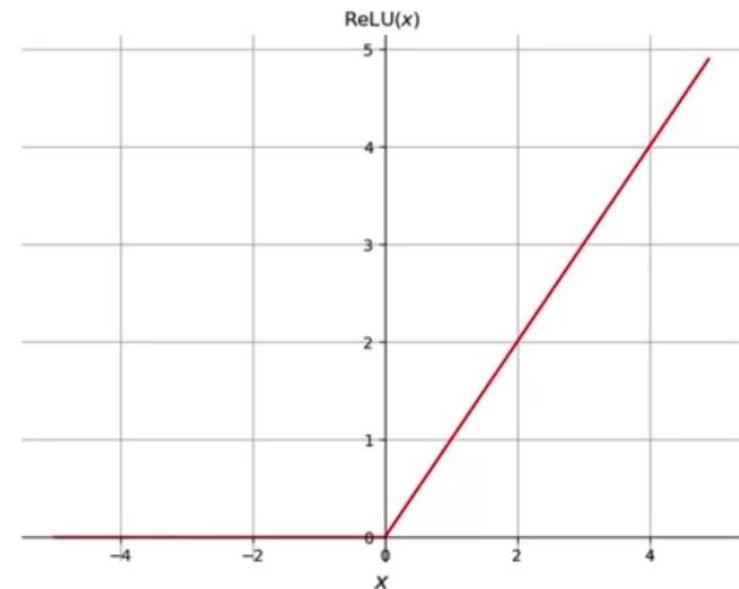
5.1
-0.3
:
-4.6
0.2

$\xrightarrow{\text{ReLU}}$

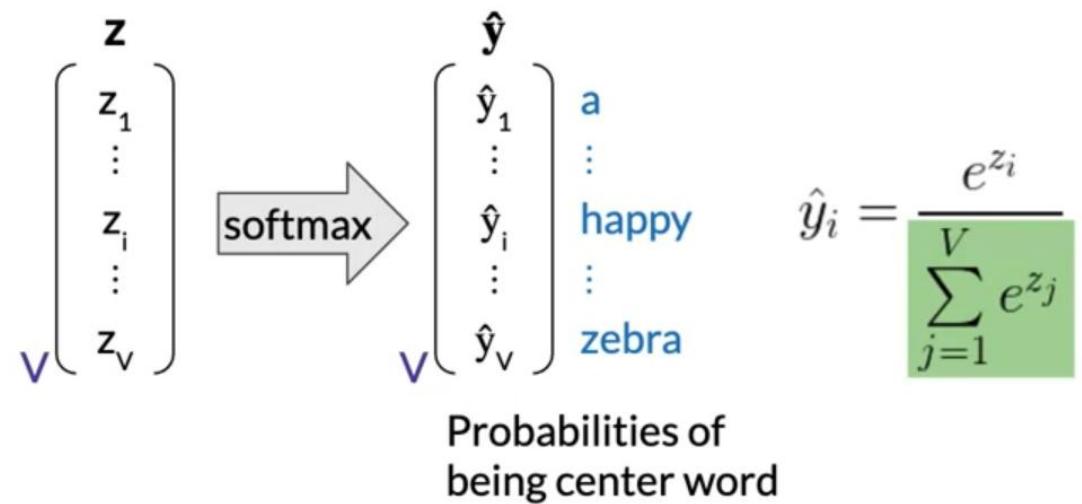
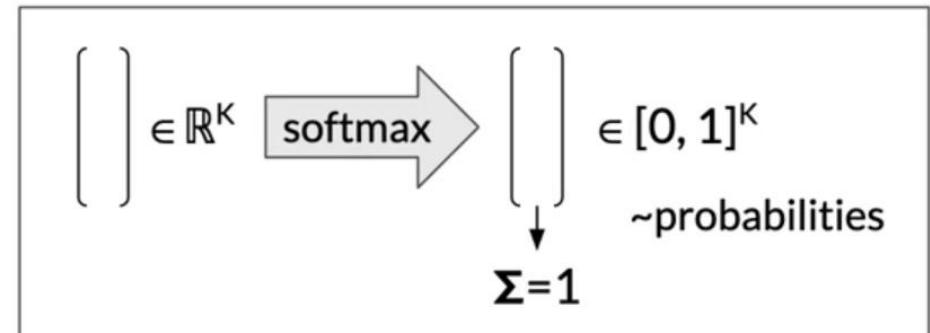
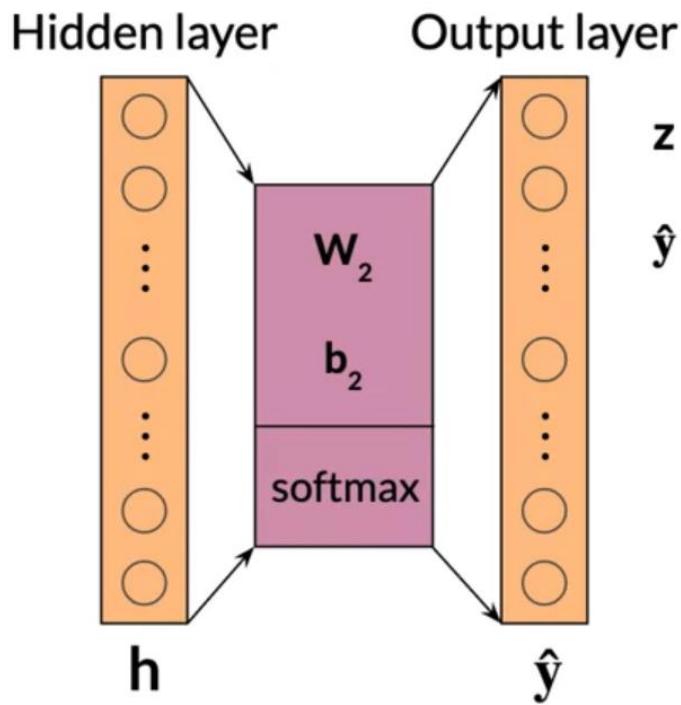
h

5.1
0
:
0
0.2

$$\text{ReLU}(x) = \max(0, x)$$

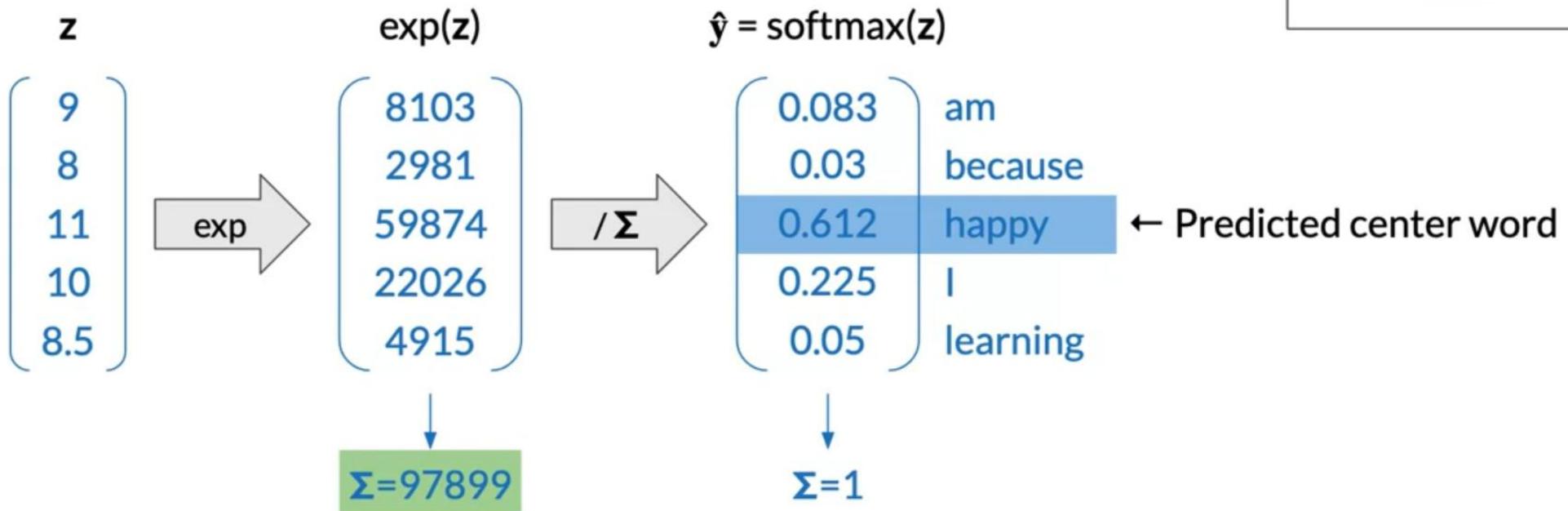


Softmax

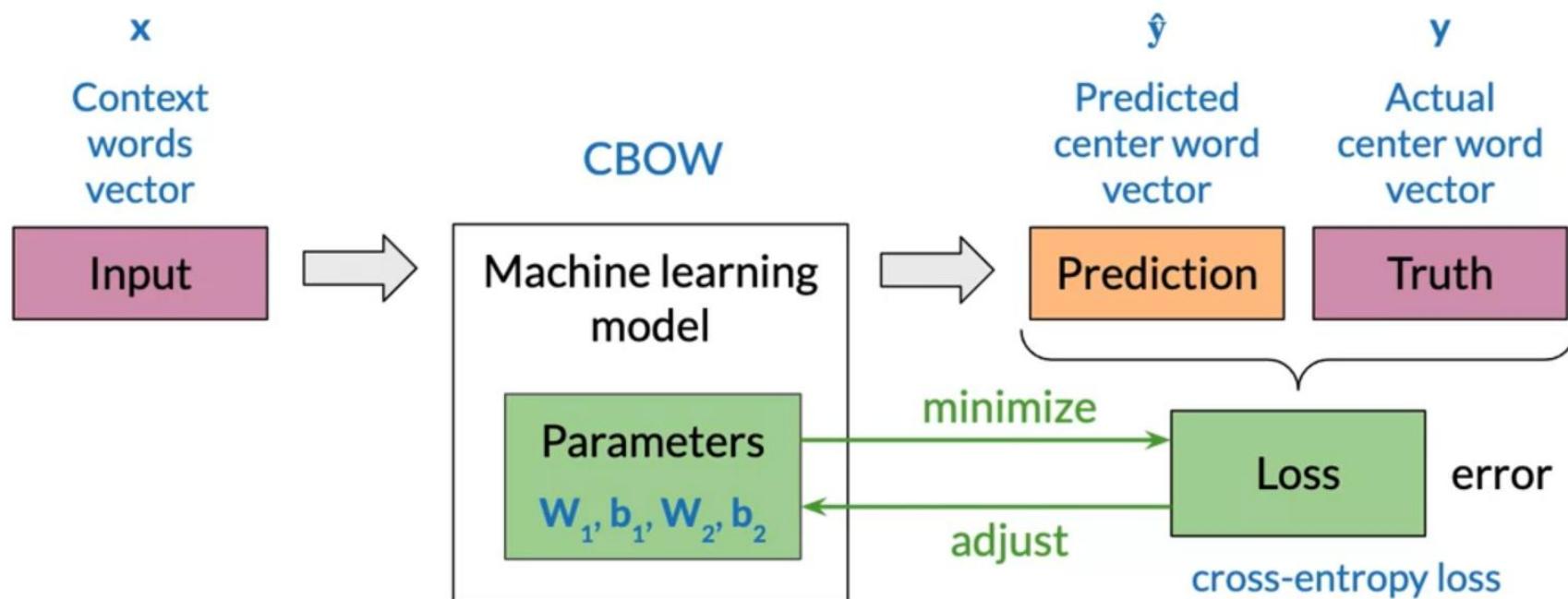


Softmax: example

$$\hat{y}_i = \frac{e^{z_i}}{\sum_{j=1}^V e^{z_j}}$$



Loss



Cross-entropy loss

$$J = - \sum_{k=1}^V y_k \log \hat{y}_k$$

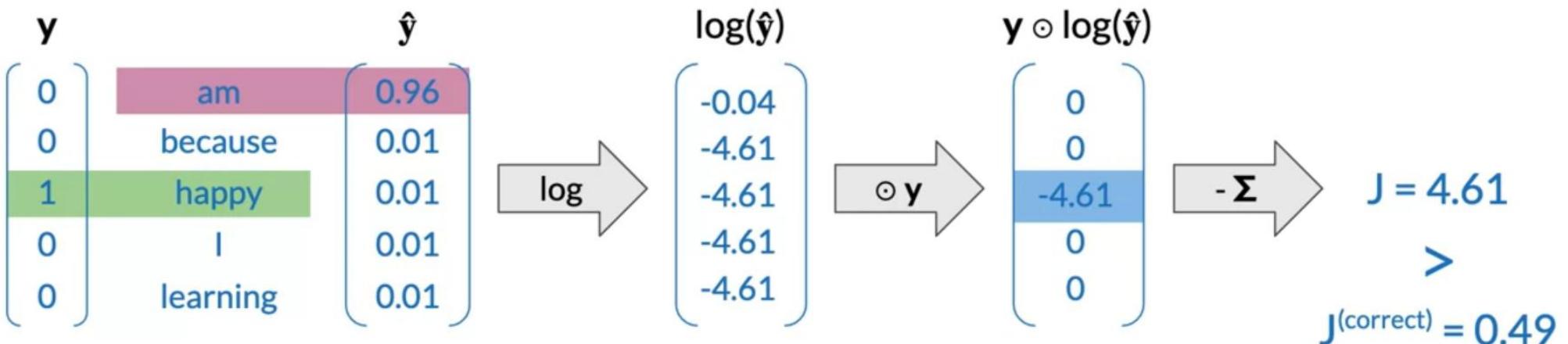
Actual $\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_v \end{bmatrix}$	Predicted $\hat{\mathbf{y}} = \begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_v \end{bmatrix}$
---	--

I am **happy** because I am learning

\mathbf{y}	$\hat{\mathbf{y}}$	$\log(\hat{\mathbf{y}})$	$\mathbf{y} \odot \log(\hat{\mathbf{y}})$	$J = -\sum$																									
<table border="1"><tr><td>0</td><td>am</td></tr><tr><td>0</td><td>because</td></tr><tr><td>1</td><td>happy</td></tr><tr><td>0</td><td>I</td></tr><tr><td>0</td><td>learning</td></tr></table>	0	am	0	because	1	happy	0	I	0	learning	<table border="1"><tr><td>0.083</td></tr><tr><td>0.03</td></tr><tr><td>0.611</td></tr><tr><td>0.225</td></tr><tr><td>0.05</td></tr></table>	0.083	0.03	0.611	0.225	0.05	<table border="1"><tr><td>-2.49</td></tr><tr><td>-3.49</td></tr><tr><td>-0.49</td></tr><tr><td>-1.49</td></tr><tr><td>-2.49</td></tr></table>	-2.49	-3.49	-0.49	-1.49	-2.49	<table border="1"><tr><td>0</td></tr><tr><td>0</td></tr><tr><td>-0.49</td></tr><tr><td>0</td></tr><tr><td>0</td></tr></table>	0	0	-0.49	0	0	0.49
0	am																												
0	because																												
1	happy																												
0	I																												
0	learning																												
0.083																													
0.03																													
0.611																													
0.225																													
0.05																													
-2.49																													
-3.49																													
-0.49																													
-1.49																													
-2.49																													
0																													
0																													
-0.49																													
0																													
0																													

Cross-entropy loss

$$J = - \sum_{k=1}^V y_k \log \hat{y}_k$$

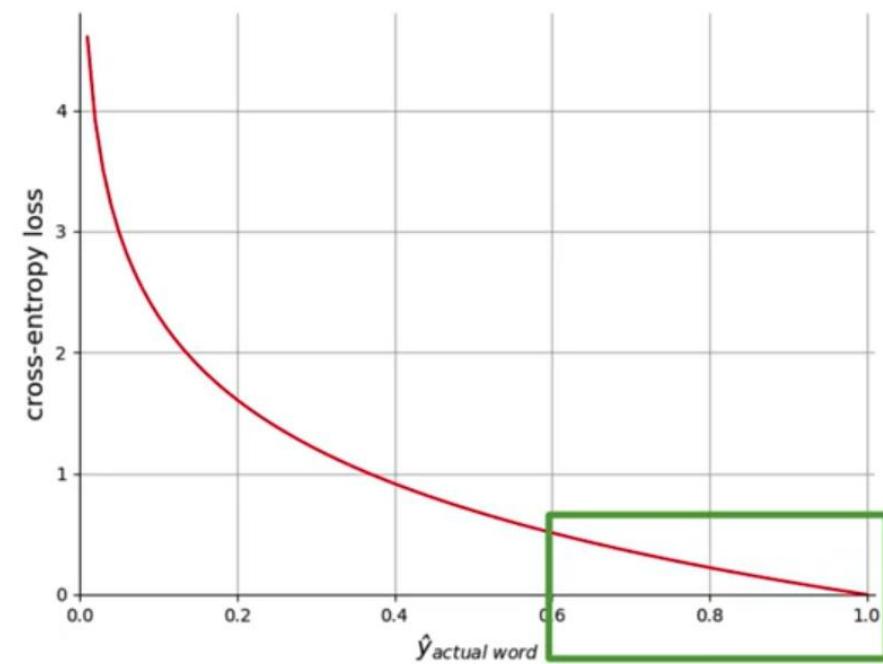


Cross-entropy loss

$$J = -\log \hat{y}_{\text{actual word}}$$

y		\hat{y}
0	am	0.96
0	because	0.01
1	happy	0.01
0	I	0.01
0	learning	0.01

$$\rightarrow J = 4.61$$

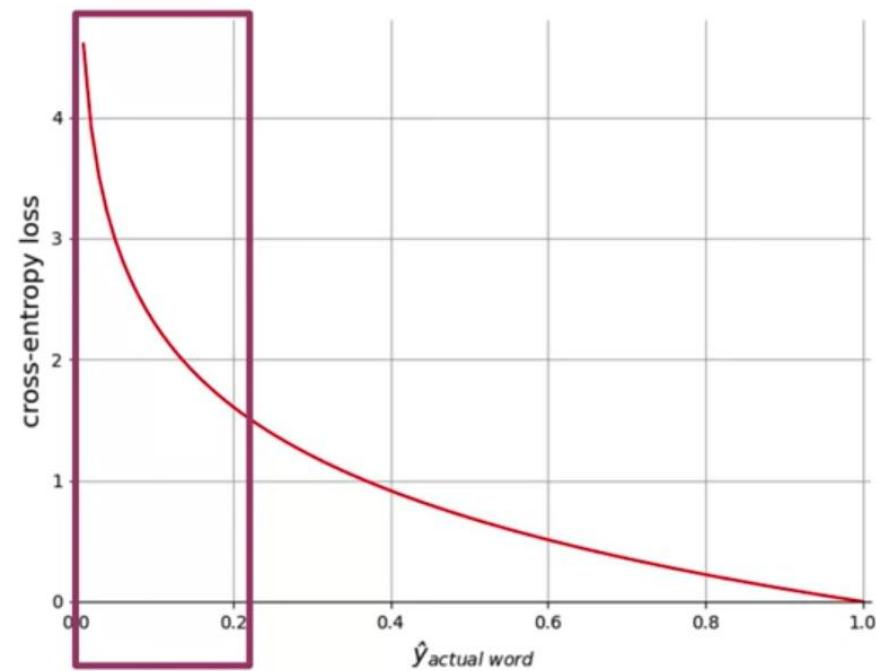


Cross-entropy loss

$$J = -\log \hat{y}_{\text{actual word}}$$

y		\hat{y}
0	am	0.96
0	because	0.01
1	happy	0.01
0	I	0.01
0	learning	0.01

$$\rightarrow J = 4.61$$



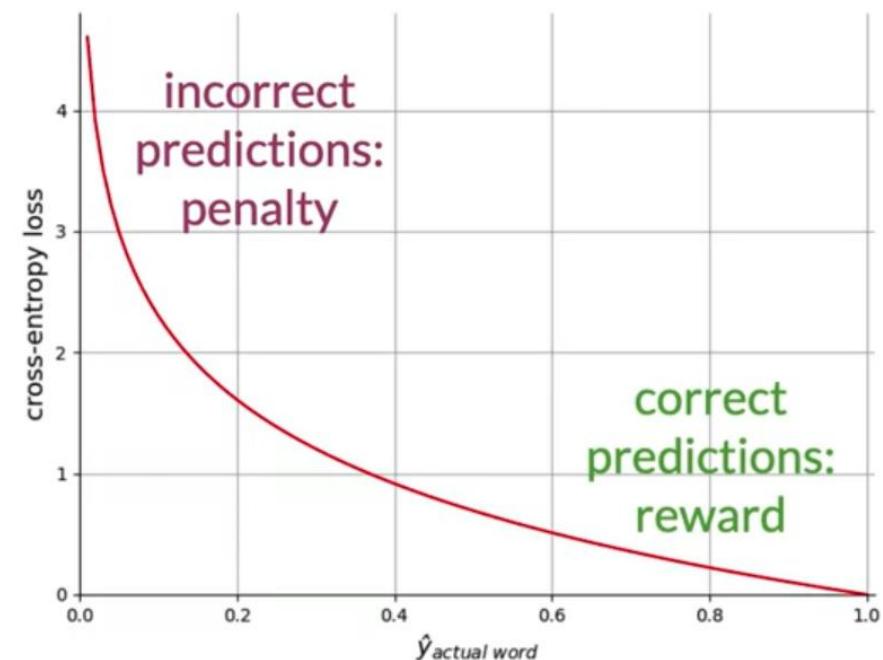
$$J = - \sum_{k=1}^V y_k \log \hat{y}_k$$

Cross-entropy loss

$$J = -\log \hat{y}_{\text{actual word}}$$

y		\hat{y}
0	am	0.96
0	because	0.01
1	happy	0.01
0	I	0.01
0	learning	0.01

$$\rightarrow J = 4.61$$



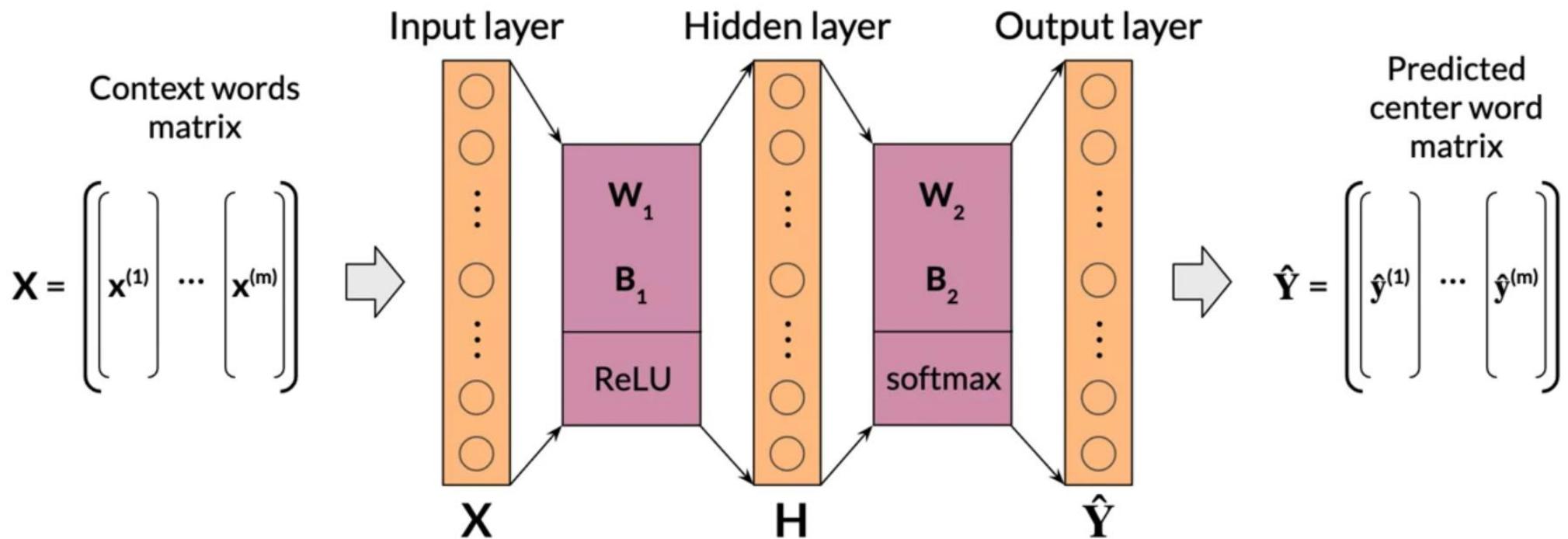
$$J = - \sum_{k=1}^V y_k \log \hat{y}_k$$

Training process

- Forward propagation
- Cost
- Backpropagation and gradient descent

Forward propagation

$$\begin{aligned} Z_1 &= \mathbf{W}_1 \mathbf{X} + \mathbf{B}_1 & Z_2 &= \mathbf{W}_2 \mathbf{H} + \mathbf{B}_2 \\ \mathbf{H} &= \text{ReLU}(Z_1) & \hat{\mathbf{Y}} &= \text{softmax}(Z_2) \end{aligned}$$



Cost

Cost: mean of losses

$$J_{batch} = -\frac{1}{m} \sum_{i=1}^m \sum_{j=1}^V y_j^{(i)} \log \hat{y}_j^{(i)}$$

$$J_{batch} = -\frac{1}{m} \sum_{i=1}^m J^{(i)}$$

Predicted
center word
matrix

$$\hat{\mathbf{Y}} = \begin{bmatrix} \hat{\mathbf{y}}^{(1)} & \cdots & \hat{\mathbf{y}}^{(m)} \end{bmatrix}$$

Actual center
word matrix

$$\mathbf{Y} = \begin{bmatrix} \mathbf{y}^{(1)} & \cdots & \mathbf{y}^{(m)} \end{bmatrix}$$

$$J = - \sum_{k=1}^V y_k \log \hat{y}_k$$

Minimizing the cost

- Backpropagation: calculate partial derivatives of cost with respect to weights and biases

$$\frac{\partial J_{batch}}{\partial \mathbf{W}_1}, \frac{\partial J_{batch}}{\partial \mathbf{W}_2}, \frac{\partial J_{batch}}{\partial \mathbf{b}_1}, \frac{\partial J_{batch}}{\partial \mathbf{b}_2}$$

- Gradient descent: update weights and biases

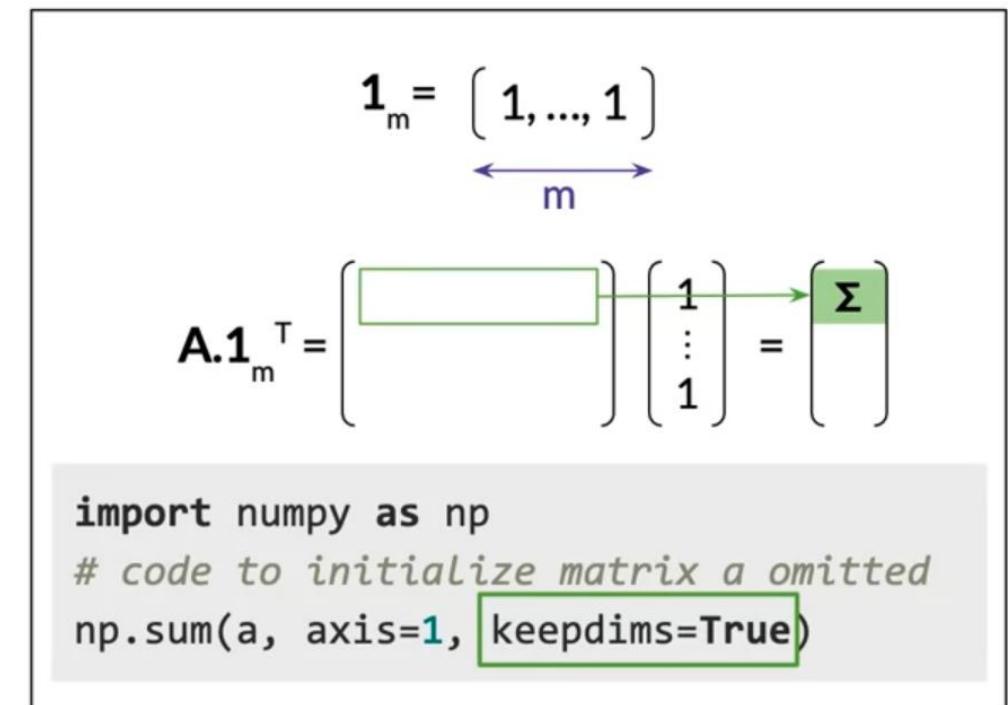
Backpropagation

$$\frac{\partial J_{batch}}{\partial \mathbf{W}_1} = \frac{1}{m} \text{ReLU} \left(\mathbf{W}_2^\top (\hat{\mathbf{Y}} - \mathbf{Y}) \right) \mathbf{X}^\top$$

$$\frac{\partial J_{batch}}{\partial \mathbf{W}_2} = \frac{1}{m} (\hat{\mathbf{Y}} - \mathbf{Y}) \mathbf{H}^\top$$

$$\frac{\partial J_{batch}}{\partial \mathbf{b}_1} = \frac{1}{m} \text{ReLU} \left(\mathbf{W}_2^\top (\hat{\mathbf{Y}} - \mathbf{Y}) \right) \mathbf{1}_m^\top$$

$$\frac{\partial J_{batch}}{\partial \mathbf{b}_2} = \frac{1}{m} (\hat{\mathbf{Y}} - \mathbf{Y}) \mathbf{1}_m^\top$$



Gradient descent

Hyperparameter: learning rate α

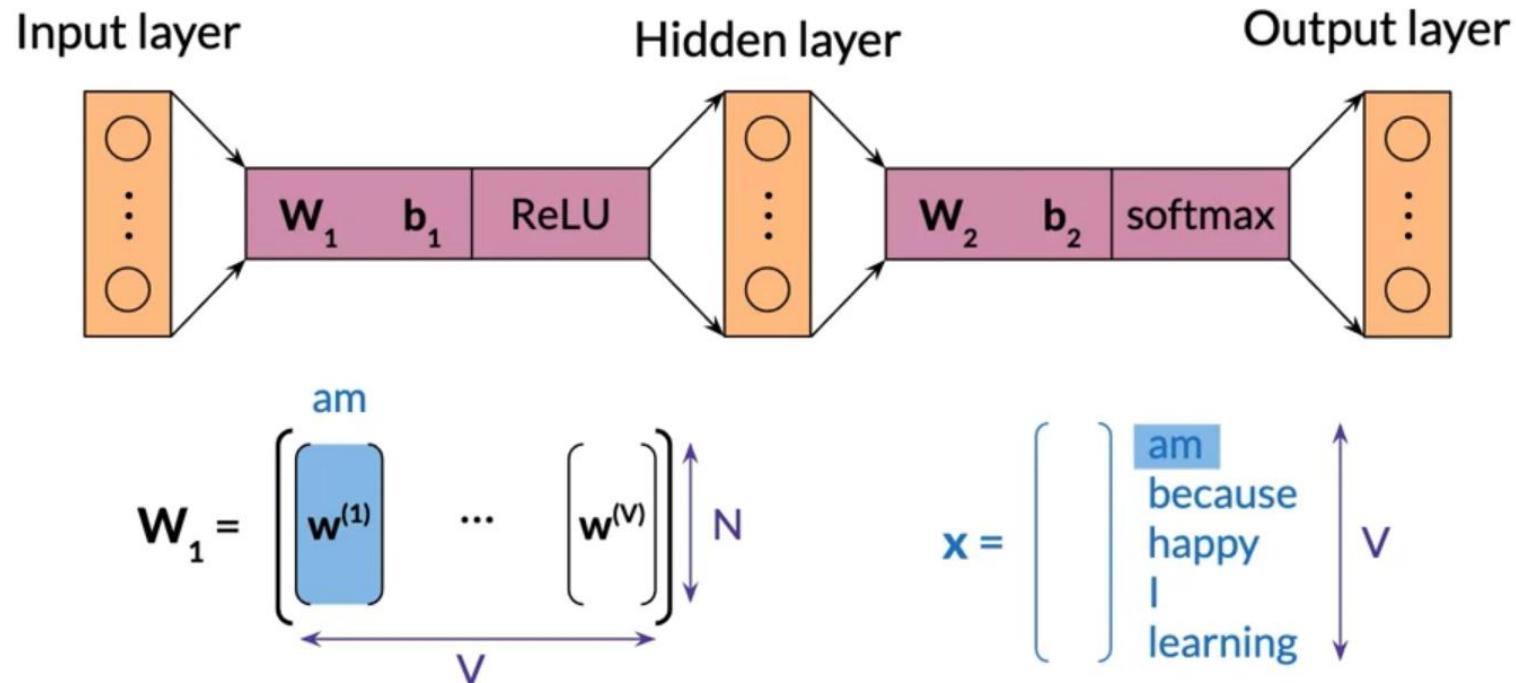
$$\mathbf{W}_1 := \mathbf{W}_1 - \alpha \frac{\partial J_{batch}}{\partial \mathbf{W}_1}$$

$$\mathbf{W}_2 := \mathbf{W}_2 - \alpha \frac{\partial J_{batch}}{\partial \mathbf{W}_2}$$

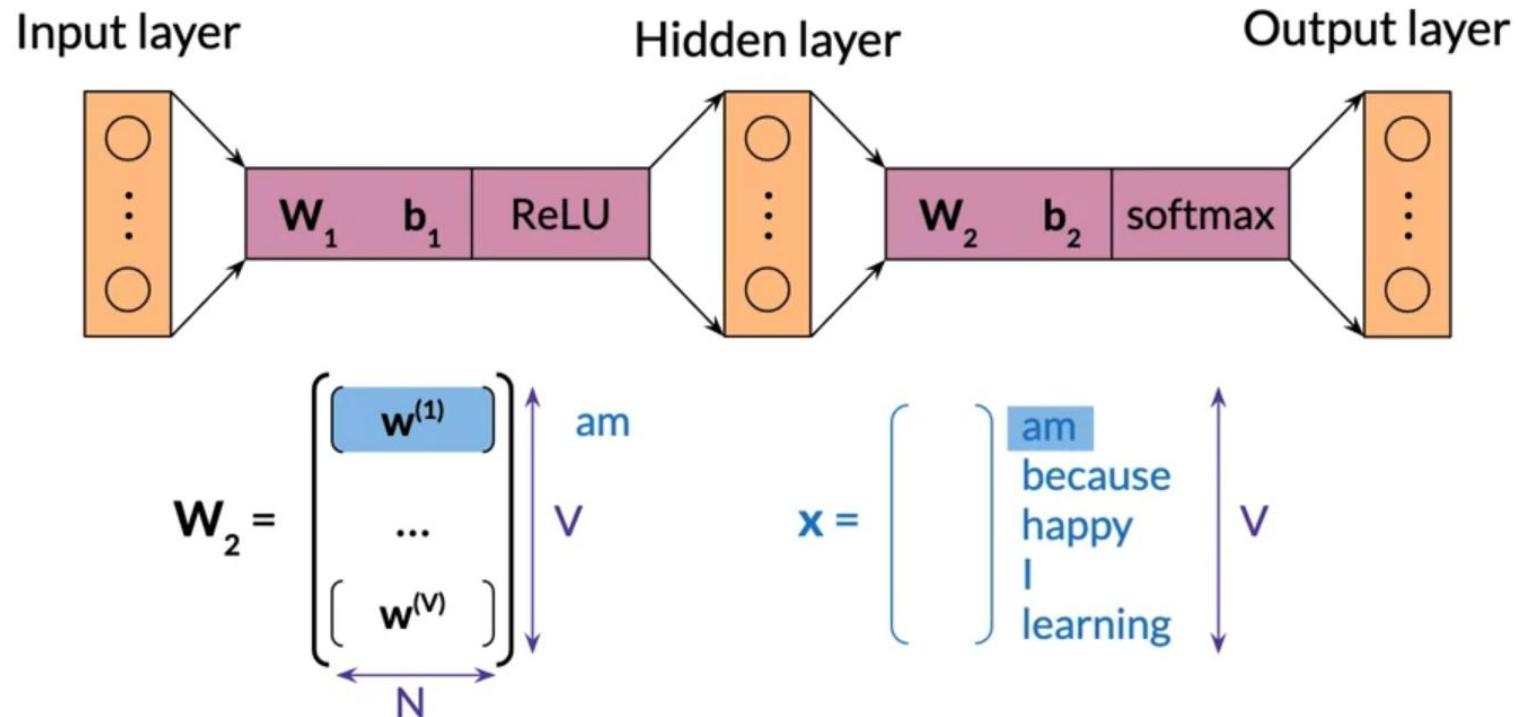
$$\mathbf{b}_1 := \mathbf{b}_1 - \alpha \frac{\partial J_{batch}}{\partial \mathbf{b}_1}$$

$$\mathbf{b}_2 := \mathbf{b}_2 - \alpha \frac{\partial J_{batch}}{\partial \mathbf{b}_2}$$

Extracting word embedding vectors: option 1



Extracting word embedding vectors: option 2



Extracting word embedding vectors: option 3

$$\mathbf{W}_1 = \begin{pmatrix} \mathbf{w}_1^{(1)} & \cdots & \mathbf{w}_1^{(V)} \end{pmatrix}$$

$$\mathbf{W}_2 = \begin{pmatrix} \mathbf{w}_2^{(1)} \\ \cdots \\ \mathbf{w}_2^{(V)} \end{pmatrix}$$

$$\mathbf{W}_3 = 0.5 (\mathbf{W}_1 + \mathbf{W}_2^T) =$$

$$\begin{pmatrix} \mathbf{w}_3^{(1)} & \cdots & \mathbf{w}_3^{(V)} \end{pmatrix}$$

$$\mathbf{x} = \begin{pmatrix} \text{am} \\ \text{because} \\ \text{happy} \\ \text{I} \\ \text{learning} \end{pmatrix}$$



Intrinsic evaluation

Test relationships between words

- Analogies

Semantic analogies

“France” is to “Paris” as “Italy” is to <?>

Syntactic analogies

“seen” is to “saw” as “been” is to <?>

 Ambiguity

“wolf” is to “pack” as “bee” is to <?> → swarm? colony?

Intrinsic evaluation

Test relationships between words

- Analogies

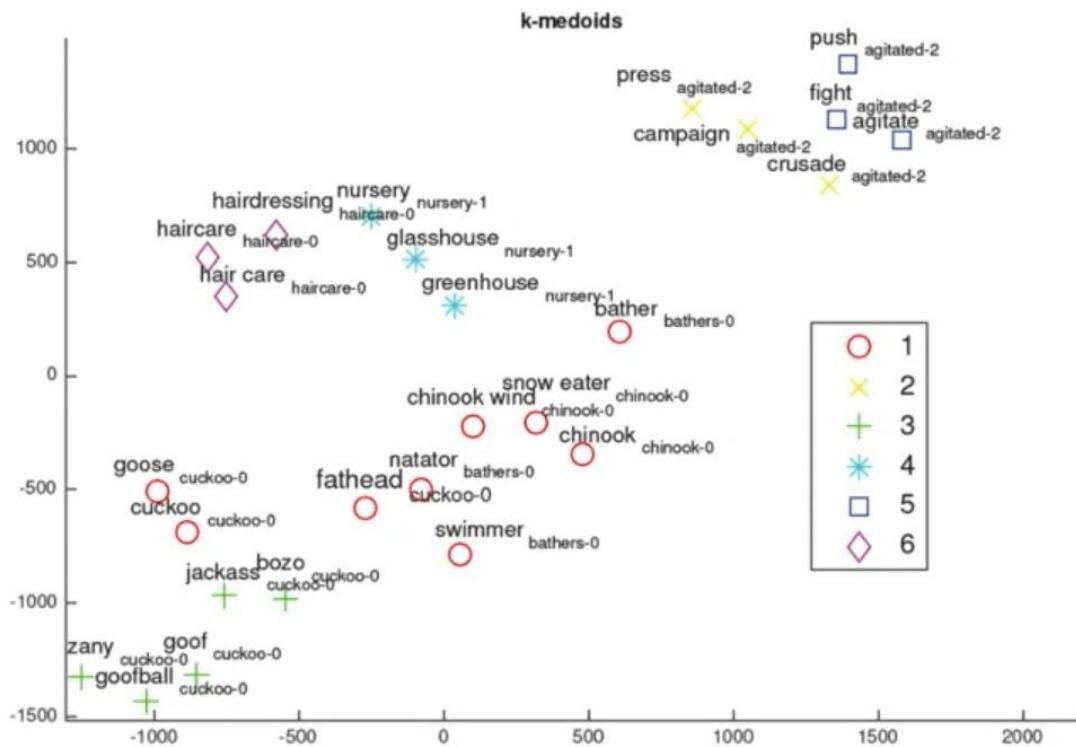
Relationship	Example 1	Example 2	Example 3
France - Paris	Italy: Rome	Japan: Tokyo	Florida: Tallahassee
big - bigger	small: larger	cold: colder	quick: quicker
Miami - Florida	Baltimore: Maryland	Dallas: Texas	Kona: Hawaii
Einstein - scientist	Messi: midfielder	Mozart: violinist	Picasso: painter
Sarkozy - France	Berlusconi: Italy	Merkel: Germany	Koizumi: Japan
copper - Cu	zinc: Zn	gold: Au	uranium: plutonium
Berlusconi - Silvio	Sarkozy: Nicolas	Putin: Medvedev	Obama: Barack
Microsoft - Windows	Google: Android	IBM: Linux	Apple: iPhone
Microsoft - Ballmer	Google: Yahoo	IBM: McNealy	Apple: Jobs
Japan - sushi	Germany: bratwurst	France: tapas	USA: pizza

Intrinsic evaluation

Test relationships between words

- Analogies
- Clustering

Source: Michael Zhai, Johnny Tan, and Jinho D. Choi. 2016. [Intrinsic and extrinsic evaluations of word embeddings](#)



Intrinsic evaluation

Test relationships between words

- Analogies
- Clustering
- Visualization



Extrinsic evaluation

Test word embeddings on external task

e.g. named entity recognition, parts-of-speech tagging

Named entity

Andrew works at deeplearning.ai

person

organization

Extrinsic evaluation

Test word embeddings on external task

e.g. named entity recognition, parts-of-speech tagging

- + Evaluates actual usefulness of embeddings
- Time-consuming
- More difficult to troubleshoot

Recap and assignment

- Data preparation
- Word representations
- Continuous bag-of-words model
- Evaluation

Going further

- Advanced language modelling and word embeddings
- NLP and machine learning libraries

Keras

```
# from keras.layers.embeddings import Embedding  
embed_layer = Embedding(10000, 400)
```

PyTorch

```
# import torch.nn as nn  
embed_layer = nn.Embedding(10000, 400)
```