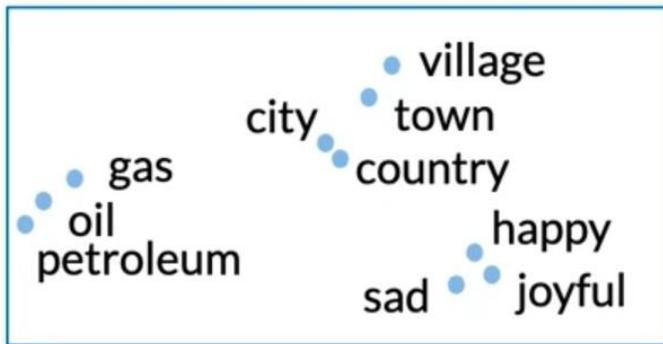


Some basic applications of word embeddings



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

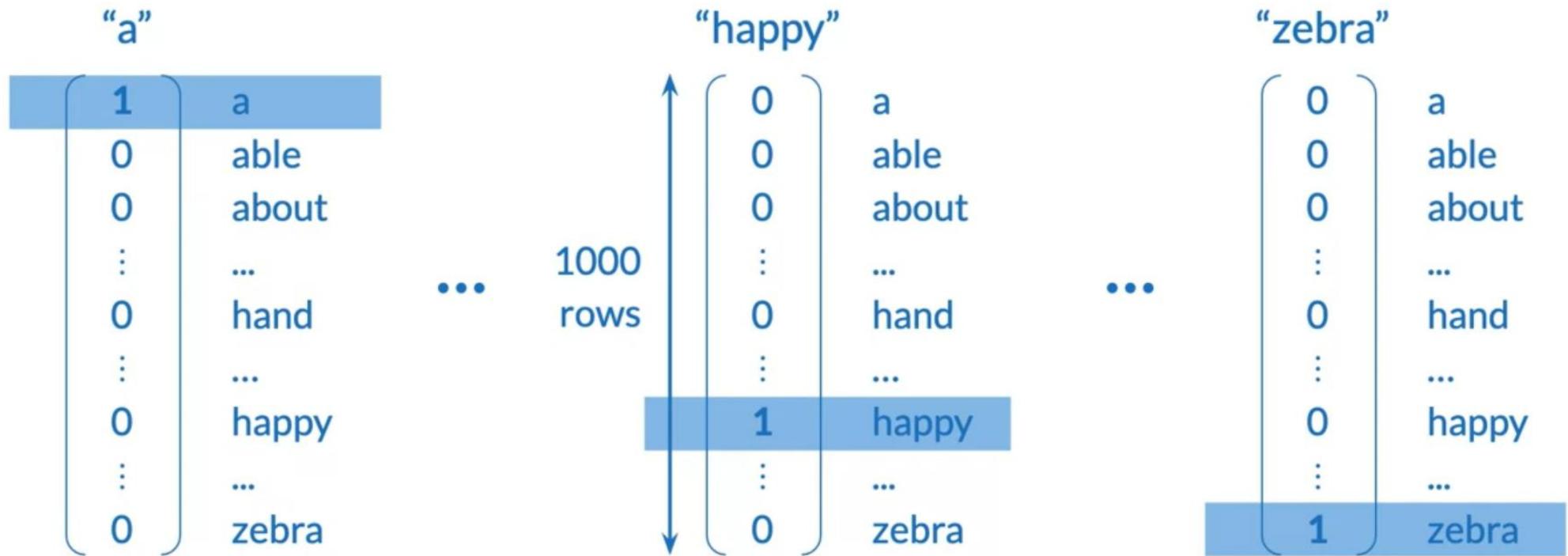
Word	Number
a	1
able	2
about	3
...	...
hand	615
...	...
happy	621
...	...
zebra	1000

Integers

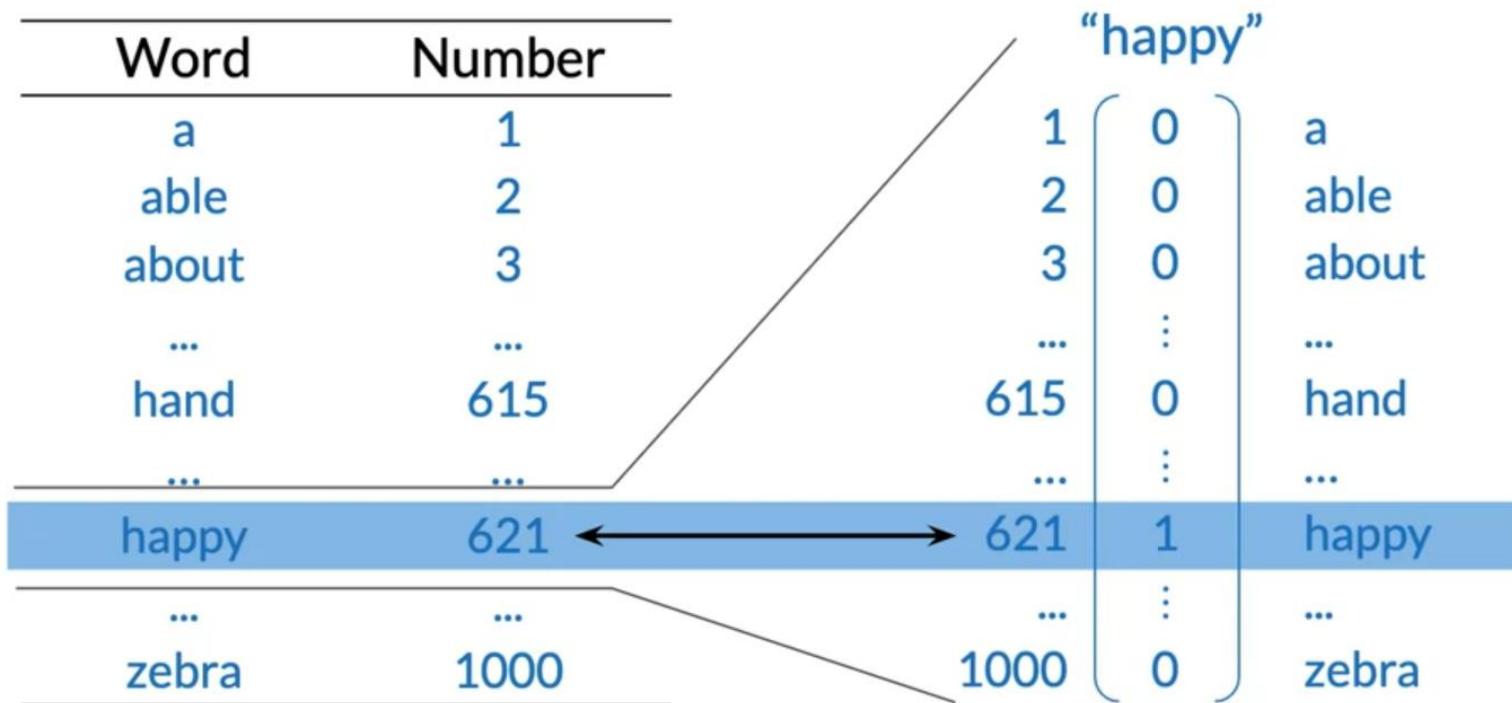
- + Simple
- Ordering: little semantic sense

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

One-hot vectors

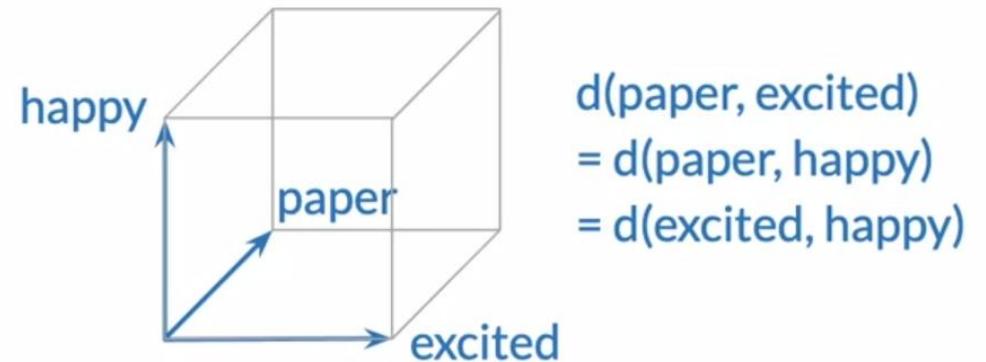
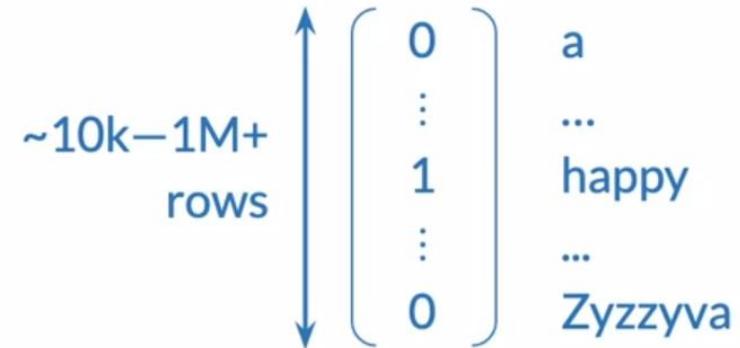


One-hot vectors

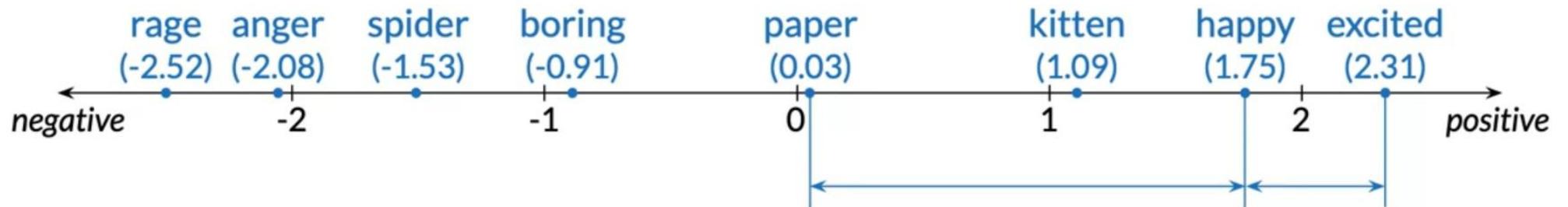


One-hot vectors

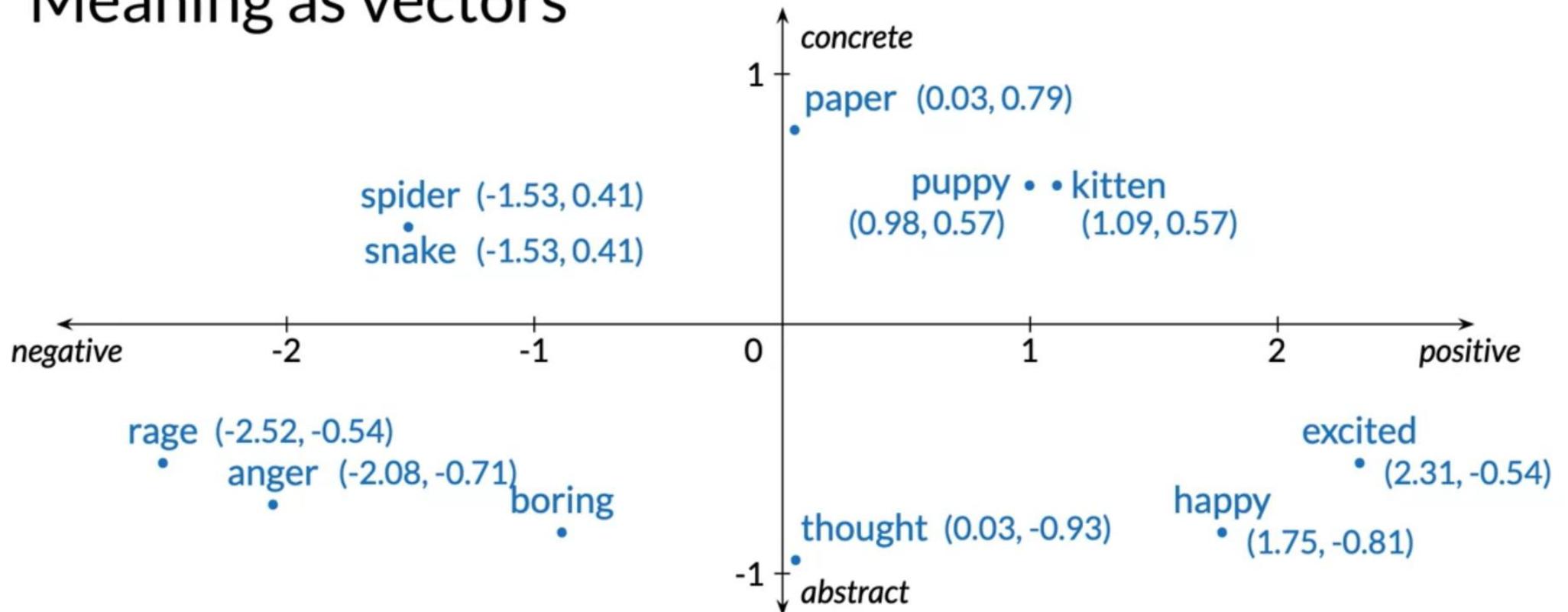
- + Simple
- + No implied ordering
- Huge vectors
- No embedded meaning



Meaning as vectors



Meaning as vectors



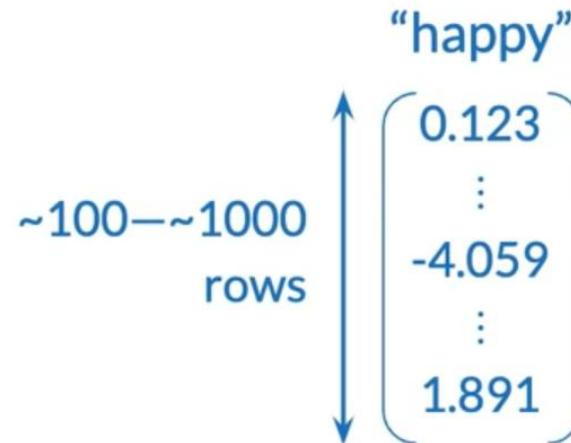
Word embedding vectors

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

forest \approx tree forest \neq ticket

- o e.g. analogies

Paris:France :: Rome:?



Terminology

integers

one-hot vectors

word vectors

word embedding vectors

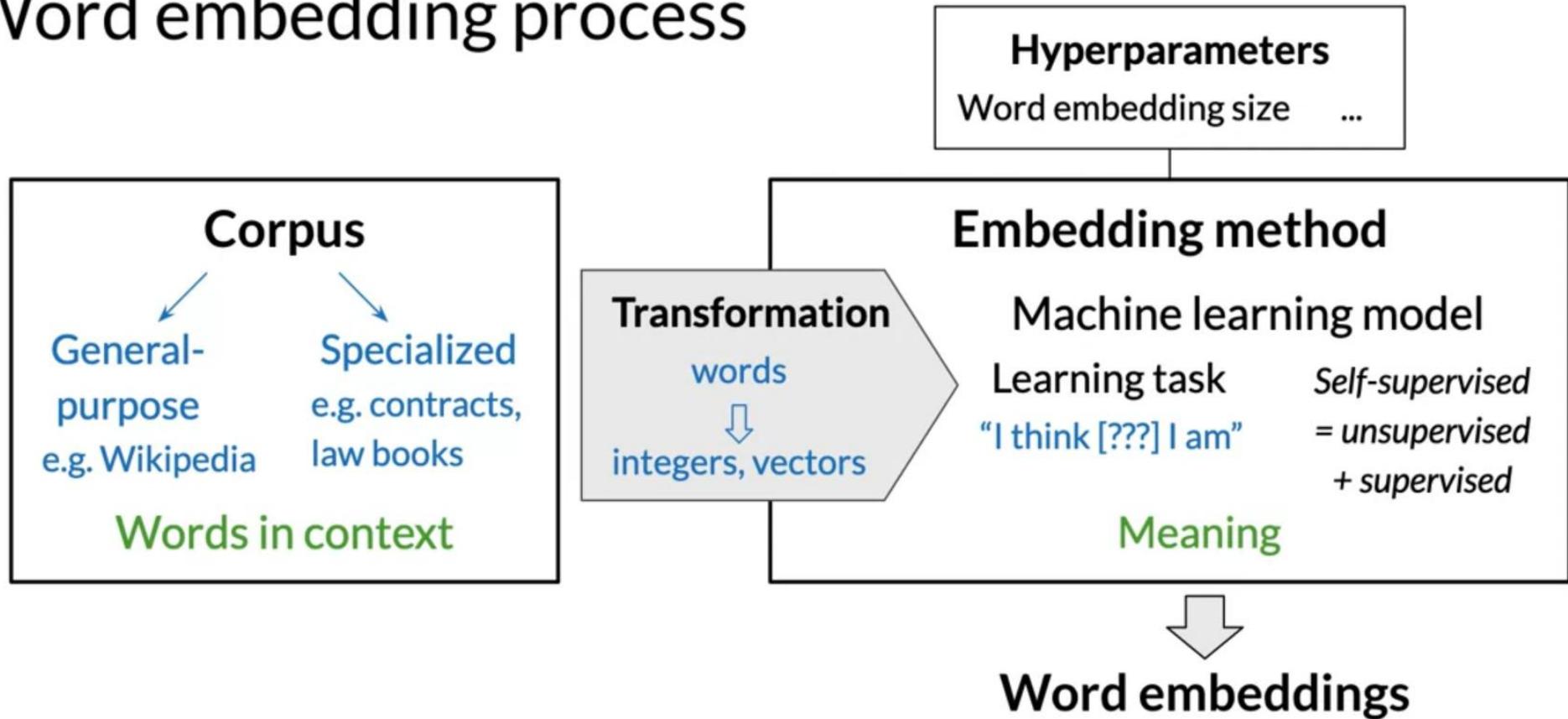
“word vectors”

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

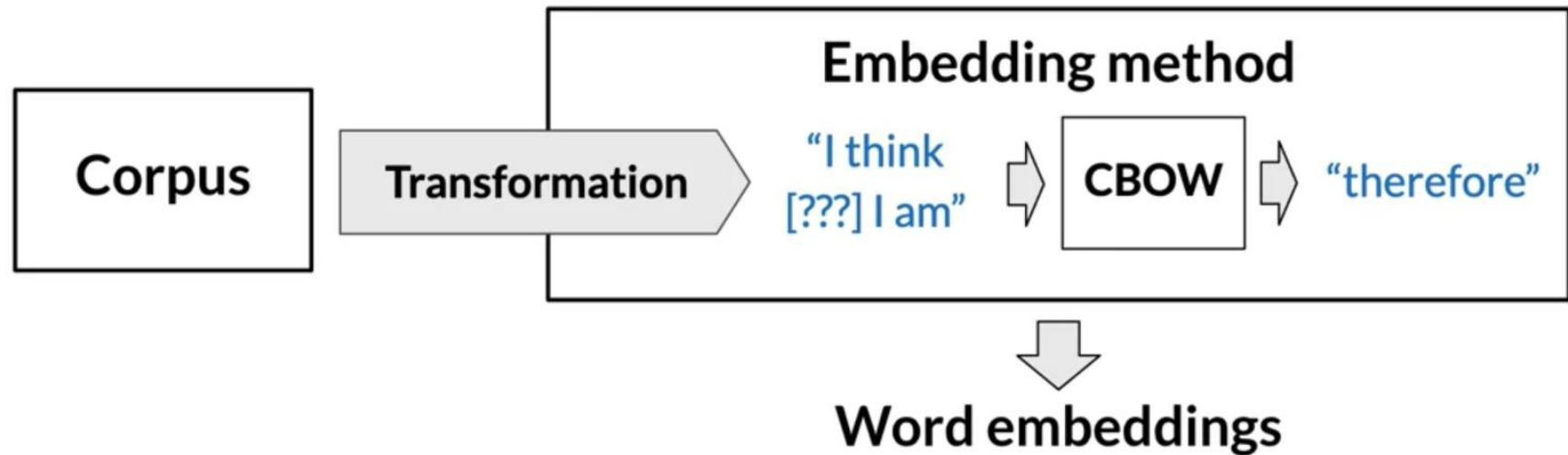
Advanced word embedding methods

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



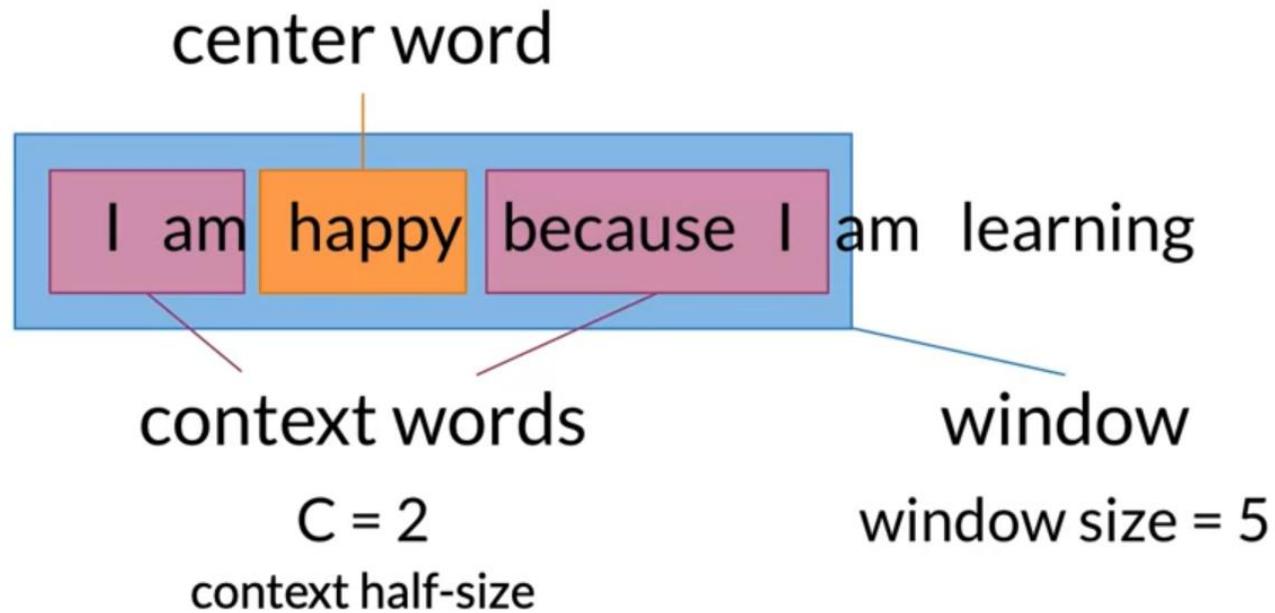
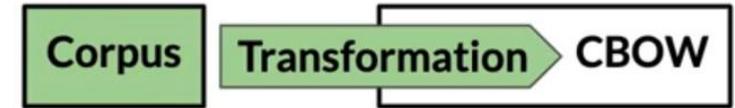
The little ? is barking



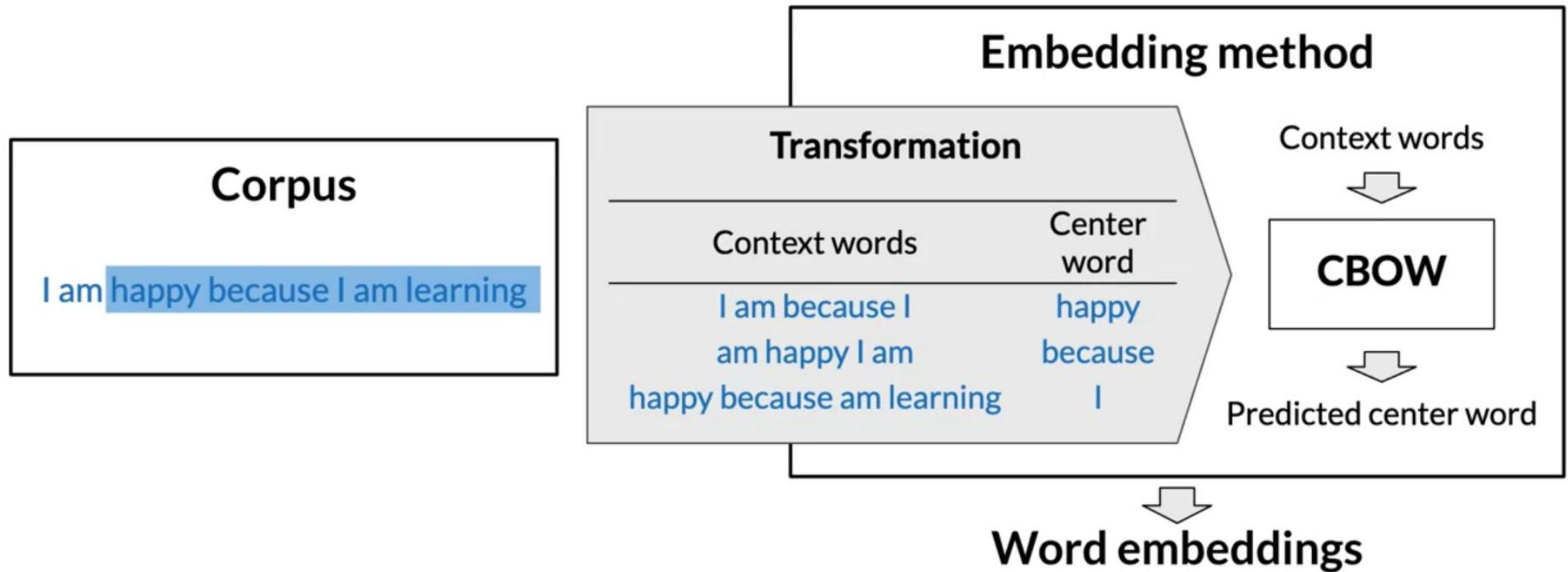
dog
puppy
hound
terrier

...

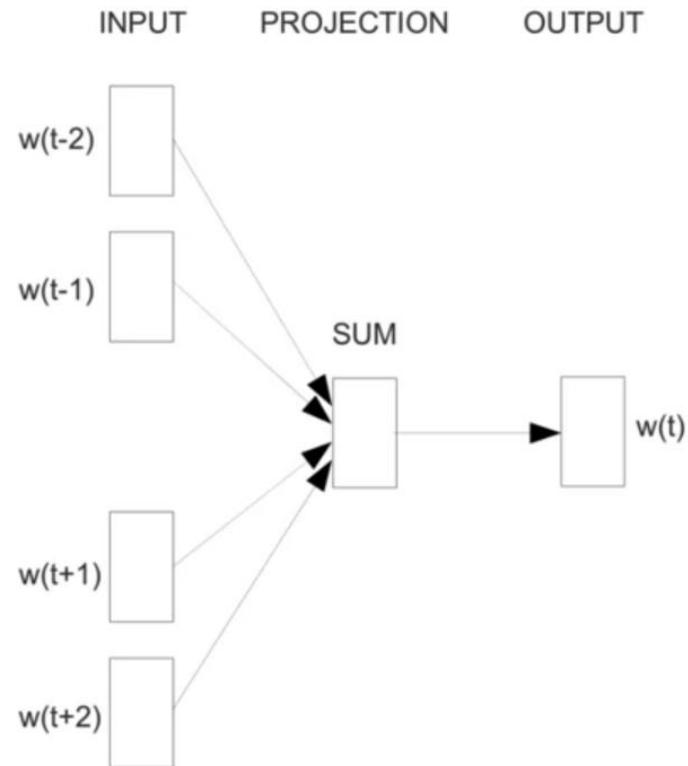
Creating a training example



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 / upper case
- Punctuation , ! . ? → . " ' « » ' " → ∅ ... !! ??? → .
- Numbers 1 2 3 5 8 → ∅ 3.14159 90210 → as is/<NUMBER>
- Special characters ∇ \$ € § ¶ ** → ∅
- 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

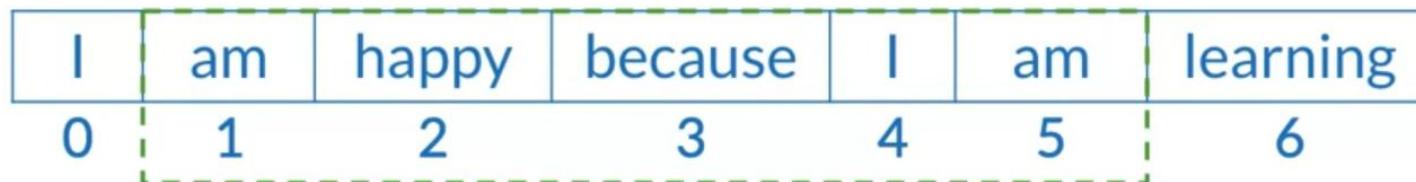
```
corpus = 'Who ❤️ "word embeddings" in 2020? I do!!!'

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	0	1	0	0	0
happy	0	0	1	0	0
I	0	0	0	1	0
learning	0	0	0	0	1

Transforming context words into vectors

Average of individual one-hot vectors

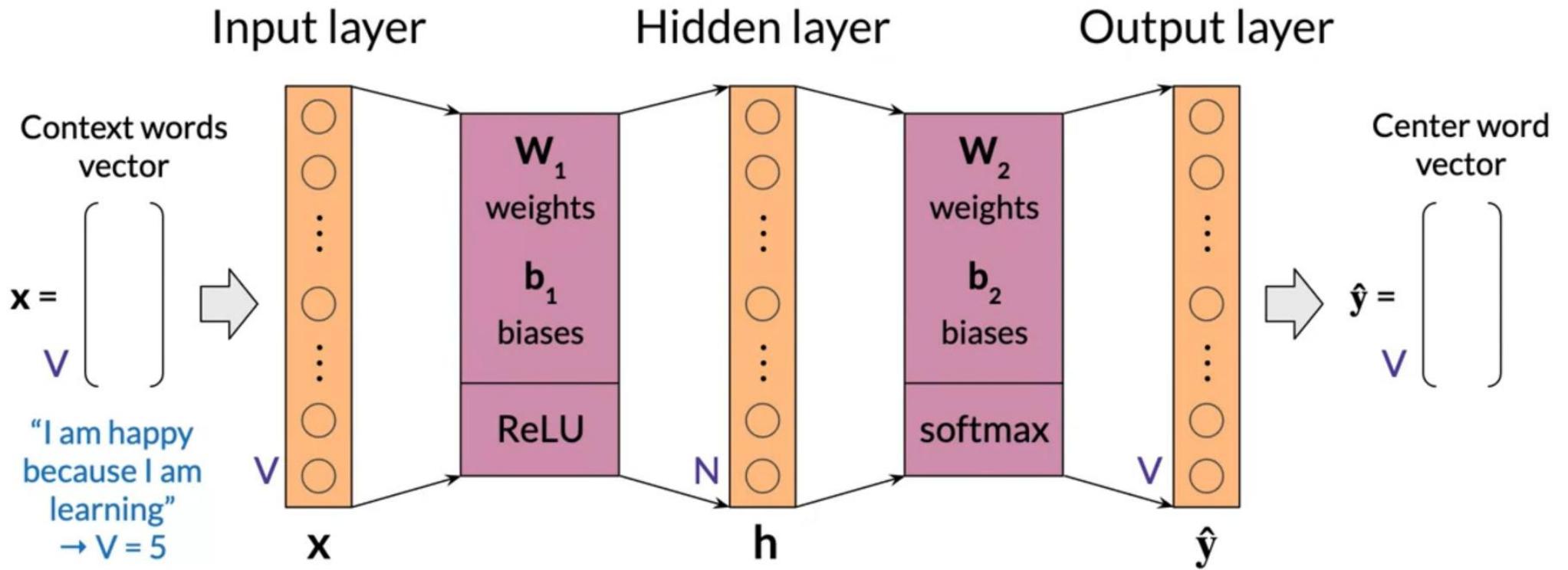
$$\left(\begin{array}{c} \text{I} \\ \text{am} \\ \text{because} \\ \text{happy} \\ \text{I} \\ \text{learning} \end{array} \begin{pmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{pmatrix} + \begin{array}{c} \text{am} \\ \text{because} \\ \text{happy} \\ \text{I} \\ \text{learning} \end{array} \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix} + \begin{array}{c} \text{because} \\ \text{am} \\ \text{happy} \\ \text{I} \\ \text{learning} \end{array} \begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{pmatrix} + \begin{array}{c} \text{I} \\ \text{am} \\ \text{because} \\ \text{happy} \\ \text{I} \\ \text{learning} \end{array} \begin{pmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{pmatrix} \right) / 4 = \begin{array}{c} \text{I am because I} \\ \begin{pmatrix} 0.25 \\ 0.25 \\ 0 \\ 0.5 \\ 0 \end{pmatrix} \end{array}$$

Final prepared training set

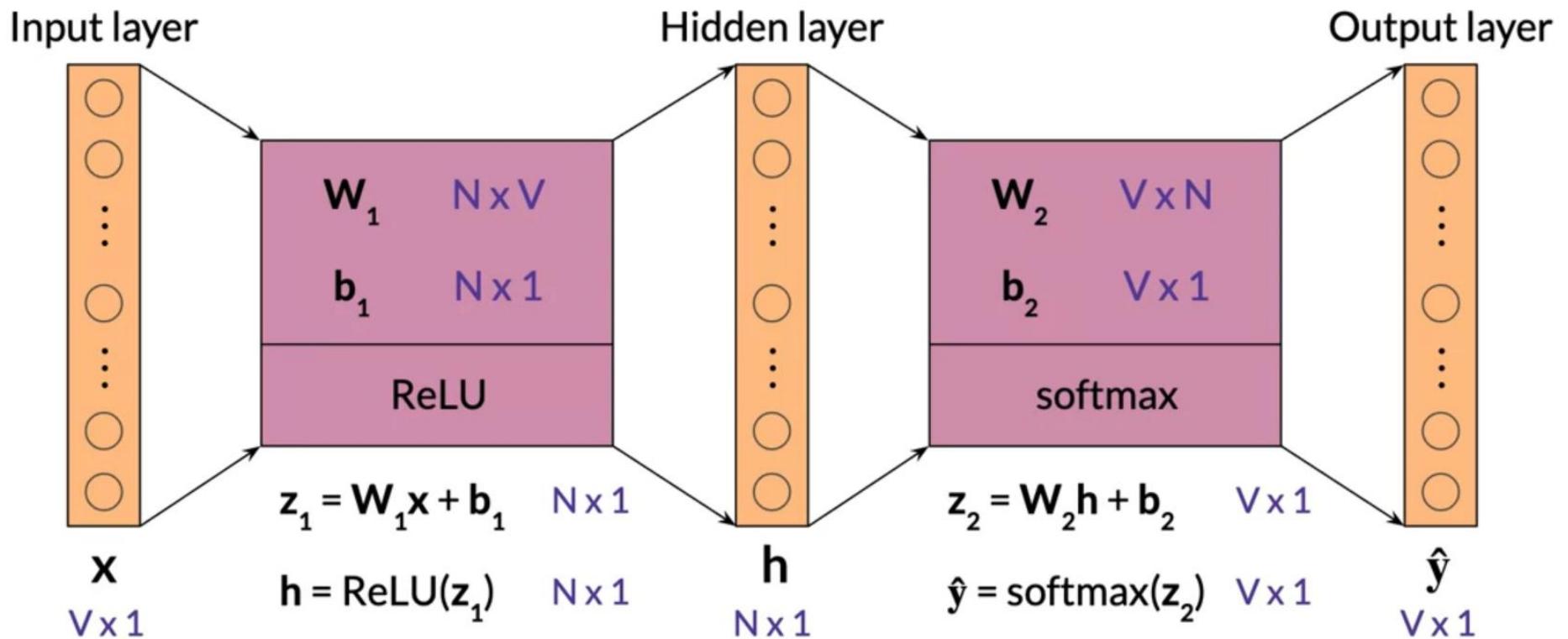
<i>Context words</i>	<i>Context words vector</i>	<i>Center word</i>	<i>Center word vector</i>
<i>I am because I</i>	$[0.25; 0.25; 0; 0.5; 0]$	<i>happy</i>	$[0; 0; 1; 0; 0]$
<i>am happy I am</i>	$[0.5; 0; 0.25; 0.25; 0]$	<i>because</i>	$[0; 1; 0; 0; 0]$
<i>happy because am learning</i>	$[0.25; 0.25; 0.25; 0; 0.25]$	<i>I</i>	$[0; 0; 0; 1; 0]$

Architecture of the CBOW model

Hyperparameters
N: Word embedding size ...



Dimensions (single input)



Dimensions (single input)

Column vectors

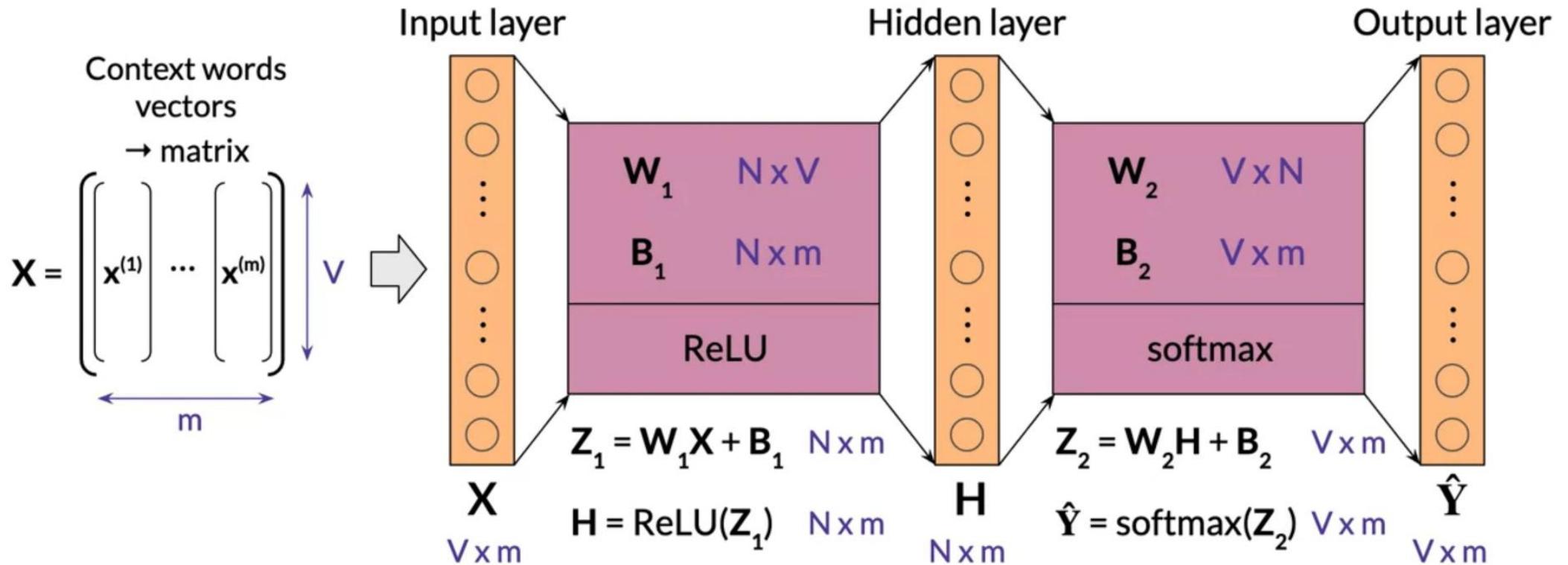
$$\mathbf{z}_1 = \mathbf{W}_1 \mathbf{x} + \mathbf{b}_1$$
$$\mathbf{z}_1 = \begin{pmatrix} \\ \end{pmatrix}_{N \times 1} \quad \mathbf{W}_1 = \begin{pmatrix} & \\ & \end{pmatrix}_{N \times V} \quad \mathbf{x} = \begin{pmatrix} \\ \end{pmatrix}_{V \times 1} \quad \mathbf{b}_1 = \begin{pmatrix} \\ \end{pmatrix}_{N \times 1}$$

Row vectors

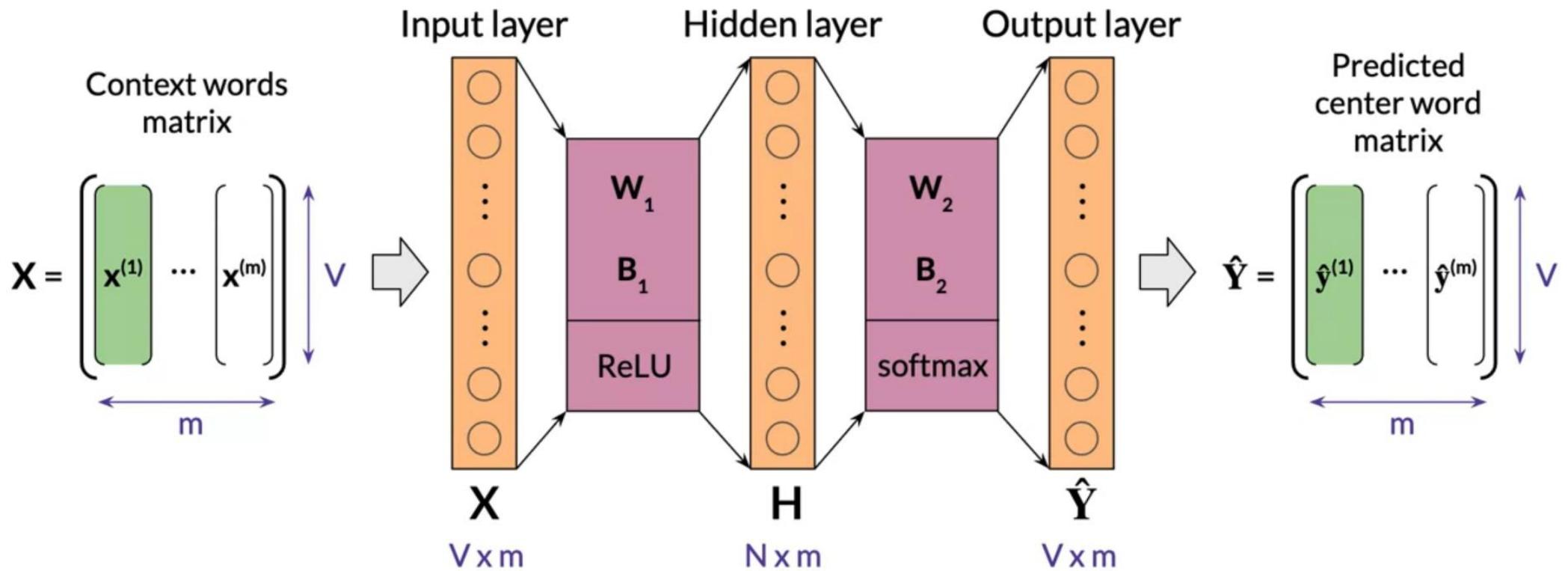
$$\mathbf{z}_1 = \mathbf{x} \mathbf{W}_1^T + \mathbf{b}_1$$
$$\mathbf{b}_1 = \begin{pmatrix} & \end{pmatrix}_{1 \times N} \quad \mathbf{W}_1 = \begin{pmatrix} & \\ & \end{pmatrix}_{N \times V} \quad \mathbf{b}_1 = \begin{pmatrix} & \end{pmatrix}_{1 \times N}$$
$$\mathbf{x} = \begin{pmatrix} & \end{pmatrix}_{1 \times V}$$

Dimensions (batch input)

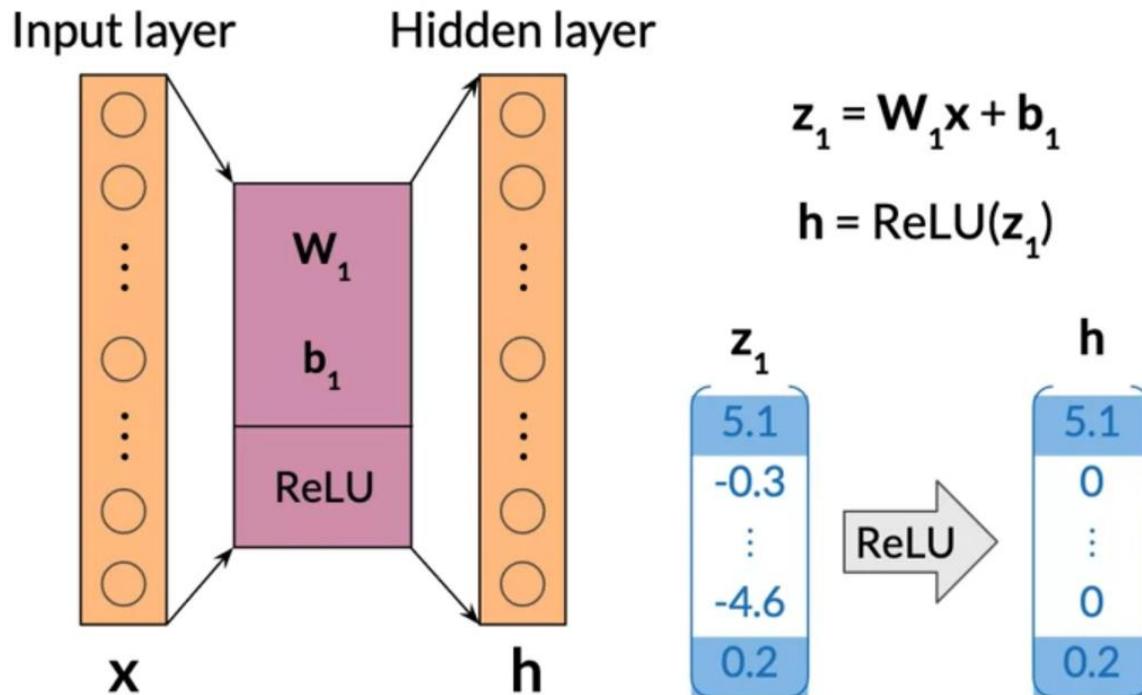
$$\begin{bmatrix} \mathbf{b}_1 \end{bmatrix} \rightarrow \mathbf{B}_1 = \begin{bmatrix} \begin{bmatrix} \mathbf{b}_1 \end{bmatrix} & \dots & \begin{bmatrix} \mathbf{b}_1 \end{bmatrix} \\ \leftarrow m \rightarrow \end{bmatrix} \begin{matrix} \updownarrow N \\ \text{broadcasting} \end{matrix}$$



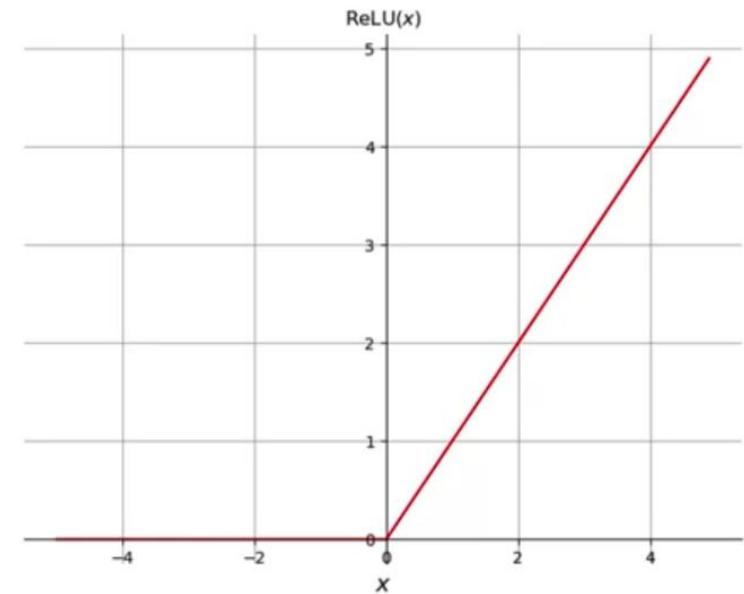
Dimensions (batch input)



Rectified Linear Unit (ReLU)



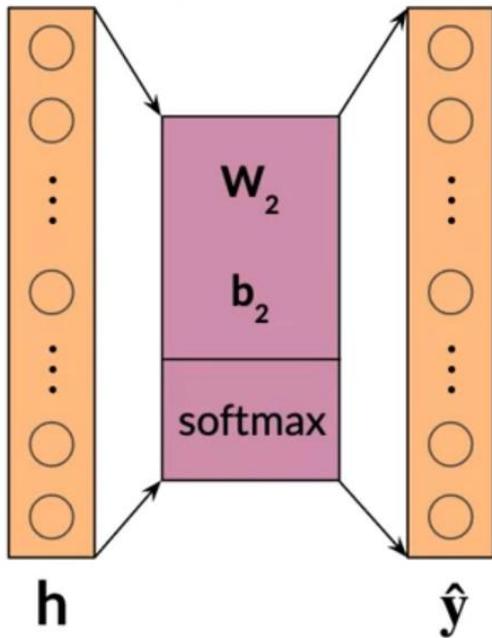
$$\text{ReLU}(x) = \max(0, x)$$



Softmax

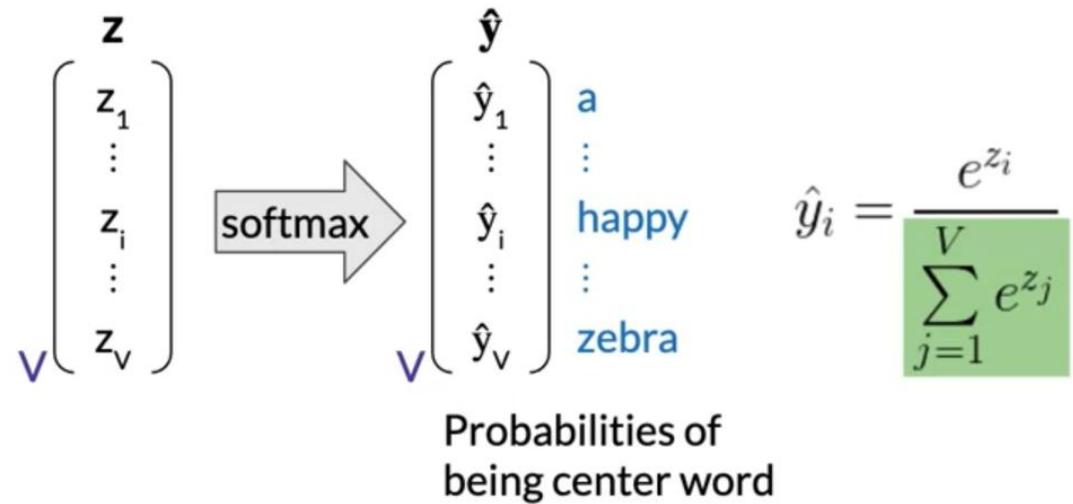
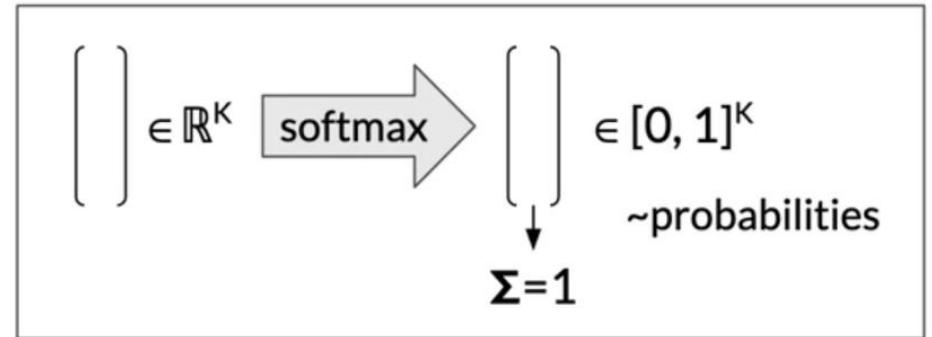
Hidden layer

Output layer

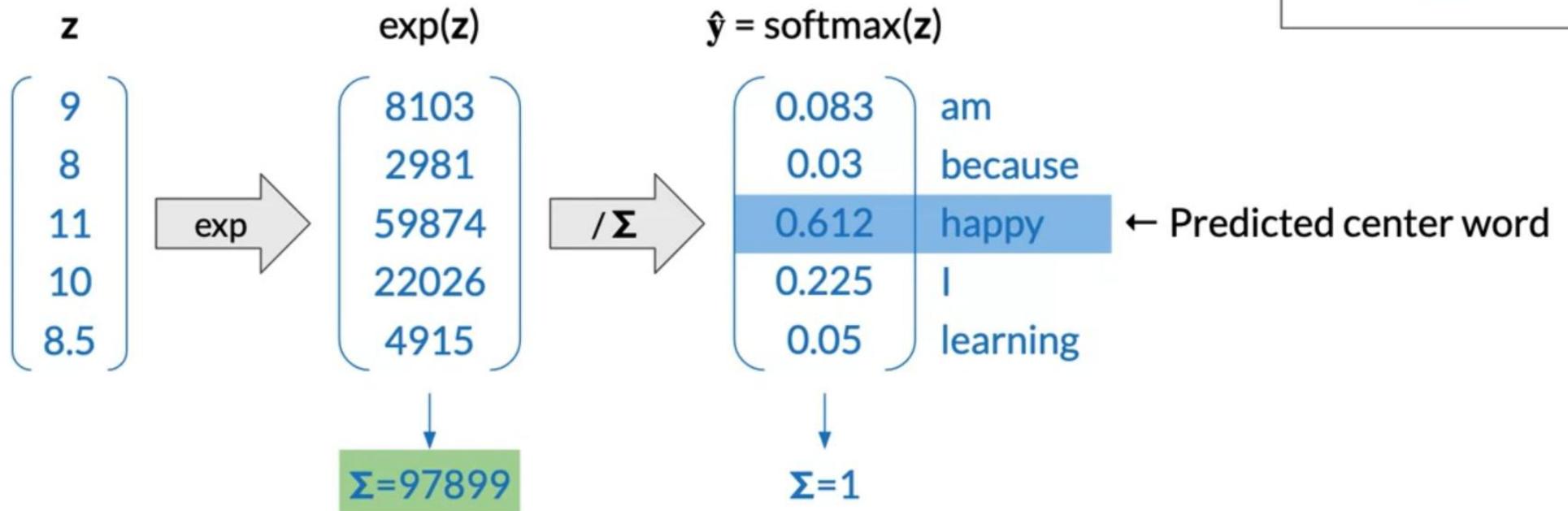


$$\mathbf{z} = \mathbf{W}_2 \mathbf{h} + \mathbf{b}_2$$

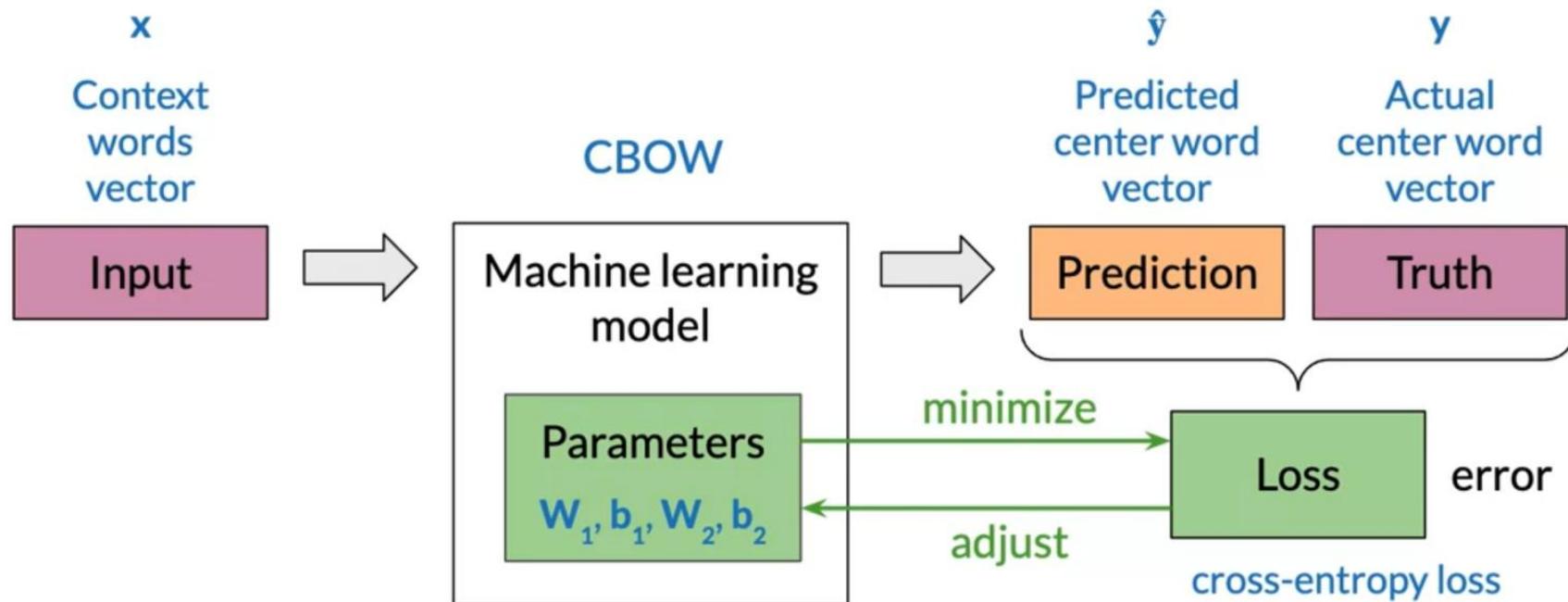
$$\hat{\mathbf{y}} = \text{softmax}(\mathbf{z})$$



Softmax: example

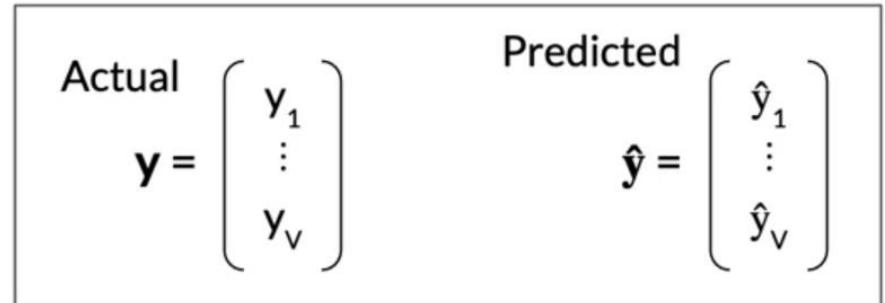


Loss

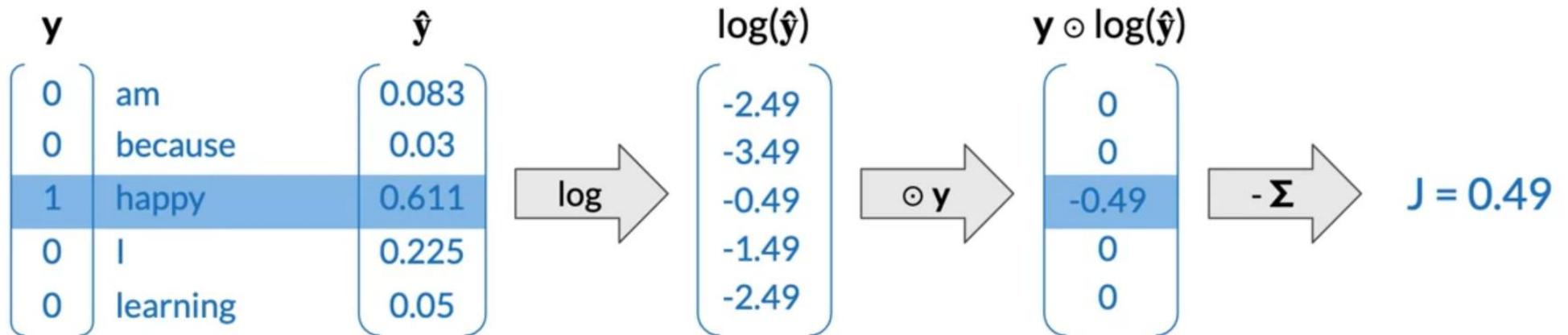


Cross-entropy loss

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

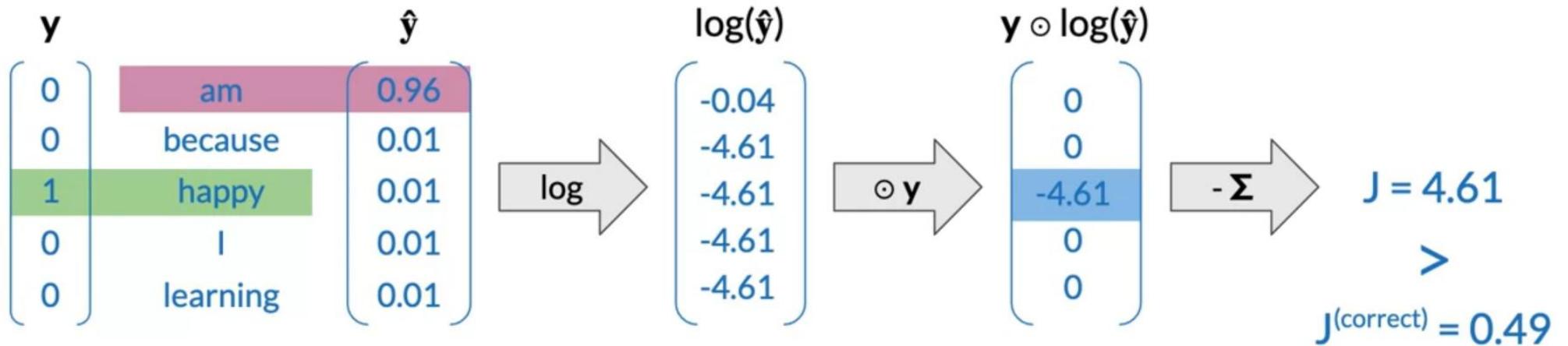


I am happy because I am learning



Cross-entropy loss

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



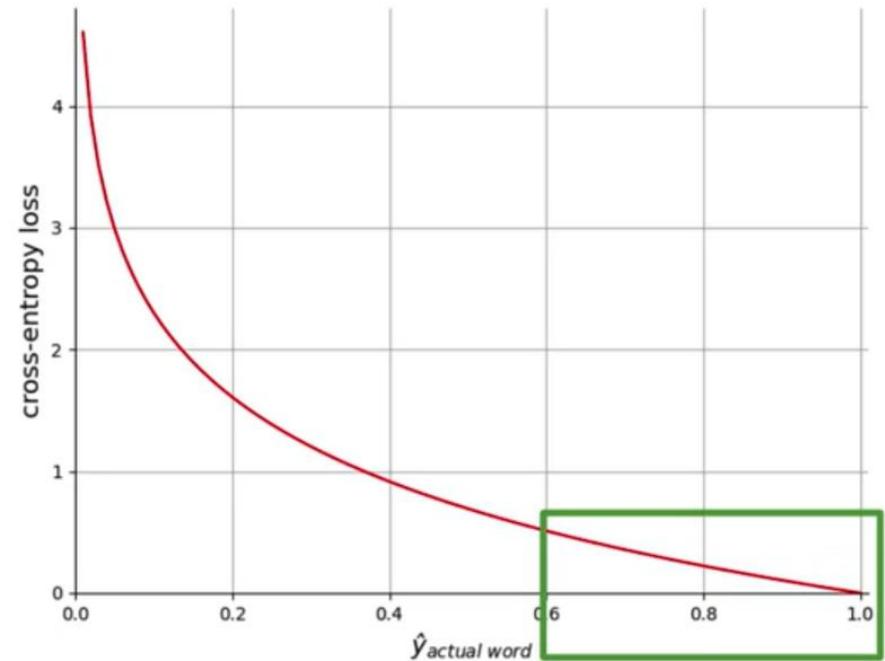
Cross-entropy loss

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

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

→ $J = 4.61$

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



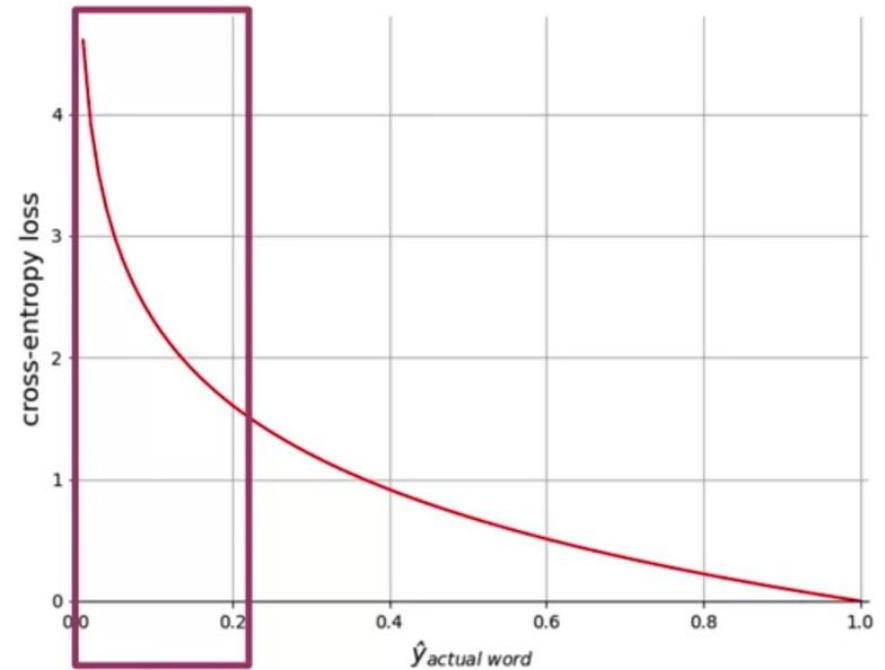
Cross-entropy loss

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

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

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

→ $J = 4.61$



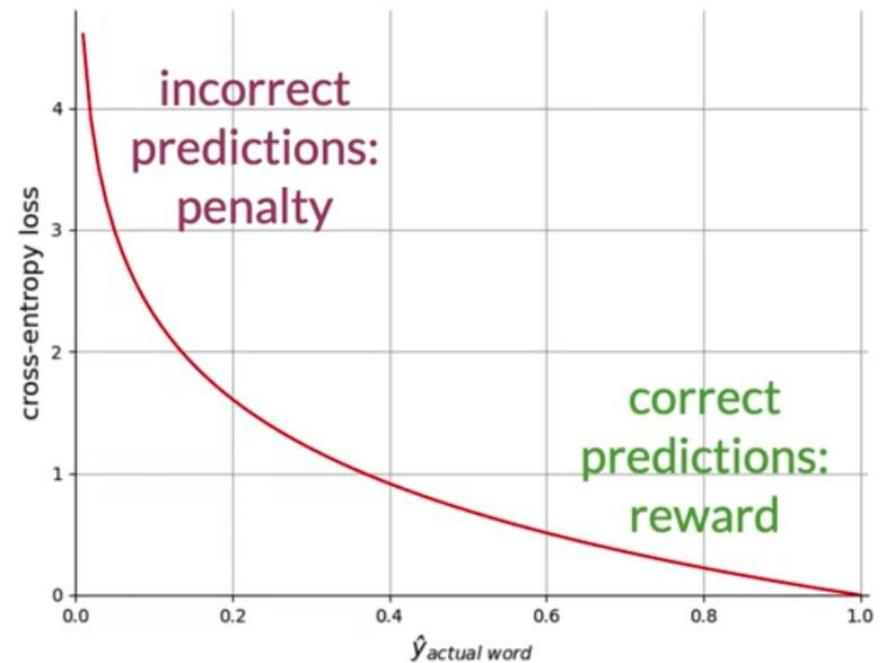
Cross-entropy loss

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

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

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

→ $J = 4.61$



Training process

- Forward propagation
- Cost
- Backpropagation and gradient descent

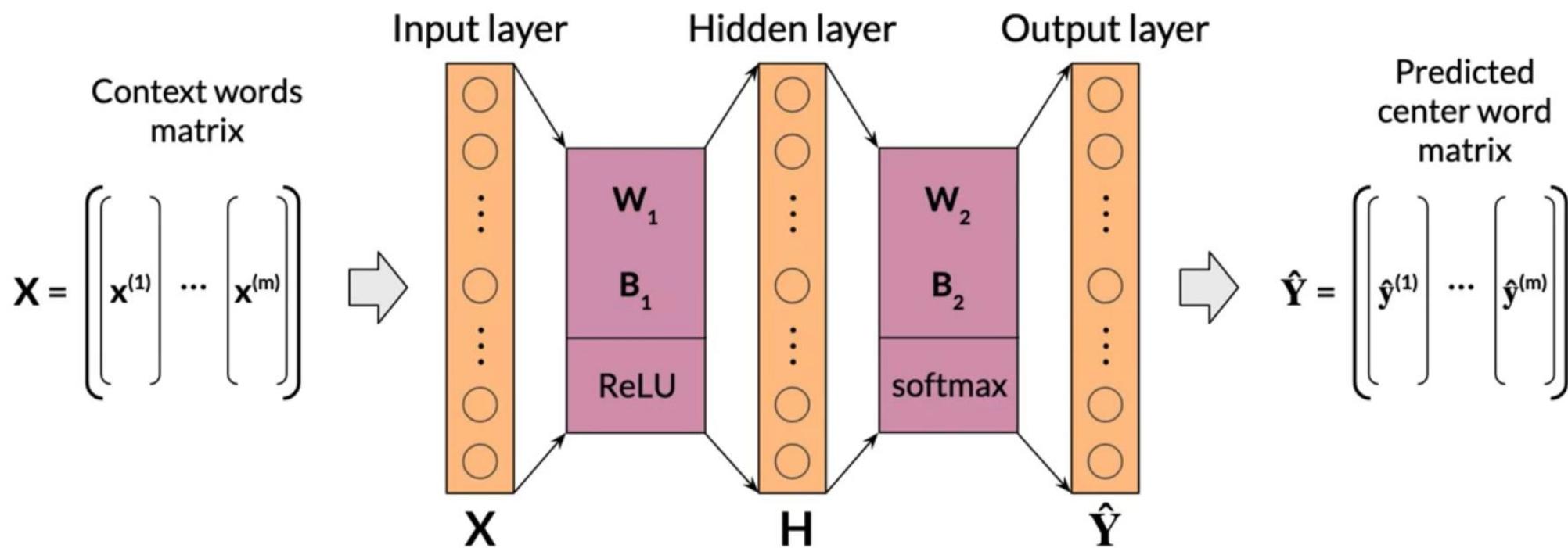
Forward propagation

$$\mathbf{Z}_1 = \mathbf{W}_1 \mathbf{X} + \mathbf{B}_1$$

$$\mathbf{Z}_2 = \mathbf{W}_2 \mathbf{H} + \mathbf{B}_2$$

$$\mathbf{H} = \text{ReLU}(\mathbf{Z}_1)$$

$$\hat{\mathbf{Y}} = \text{softmax}(\mathbf{Z}_2)$$



Cost

Cost: mean of losses

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

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

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

Predicted
center word
matrix

$$\hat{\mathbf{Y}} = \begin{bmatrix} \left[\hat{\mathbf{y}}^{(1)} \right] & \dots & \left[\hat{\mathbf{y}}^{(m)} \right] \end{bmatrix}$$

Actual center
word matrix

$$\mathbf{Y} = \begin{bmatrix} \left[\mathbf{y}^{(1)} \right] & \dots & \left[\mathbf{y}^{(m)} \right] \end{bmatrix}$$

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

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

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

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

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

$\mathbf{1}_m = \left[\underset{\substack{\longleftarrow \\ m}}{\mathbf{1}, \dots, \mathbf{1}} \right]$

$\mathbf{A} \cdot \mathbf{1}_m^\top = \begin{bmatrix} \boxed{} \\ \vdots \\ \end{bmatrix} \begin{bmatrix} \mathbf{1} \\ \vdots \\ \mathbf{1} \end{bmatrix} = \begin{bmatrix} \mathbf{\Sigma} \\ \vdots \\ \end{bmatrix}$

```
import numpy as np
# code to initialize matrix a omitted
np.sum(a, axis=1, keepdims=True)
```

Gradient descent

Hyperparameter: learning rate α

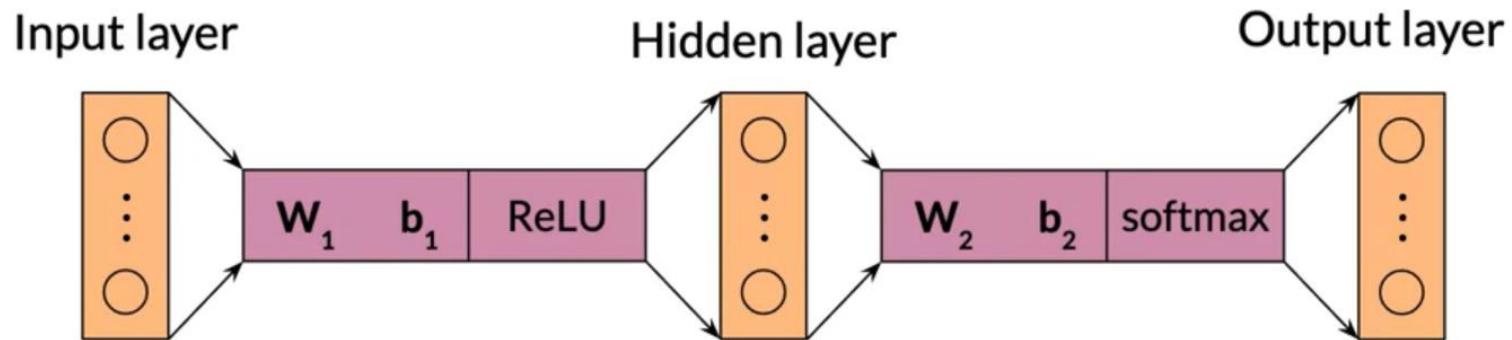
$$\mathbf{W}_1 := \mathbf{W}_1 - \alpha \frac{\partial J_{batch}}{\partial \mathbf{W}_1}$$

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

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

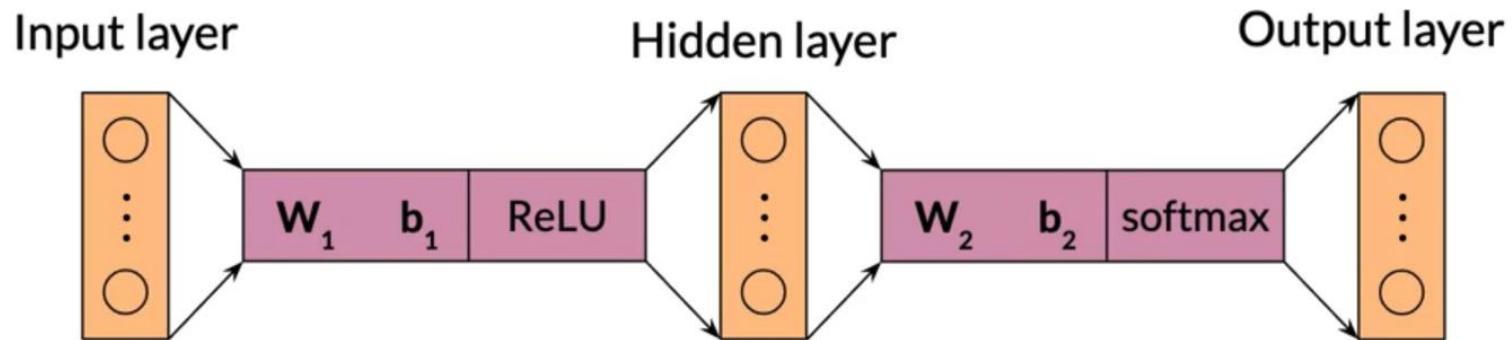
$$\mathbf{b}_2 := \mathbf{b}_2 - \alpha \frac{\partial J_{batch}}{\partial \mathbf{b}_2}$$

Extracting word embedding vectors: option 1



$$W_1 = \begin{bmatrix} \text{am} \\ \mathbf{w}^{(1)} & \dots & \mathbf{w}^{(M)} \end{bmatrix} \begin{matrix} \updownarrow N \\ \leftarrow V \end{matrix}$$
$$\mathbf{x} = \begin{bmatrix} \text{am} \\ \text{because} \\ \text{happy} \\ \text{I} \\ \text{learning} \end{bmatrix} \begin{matrix} \updownarrow V \end{matrix}$$

Extracting word embedding vectors: option 2



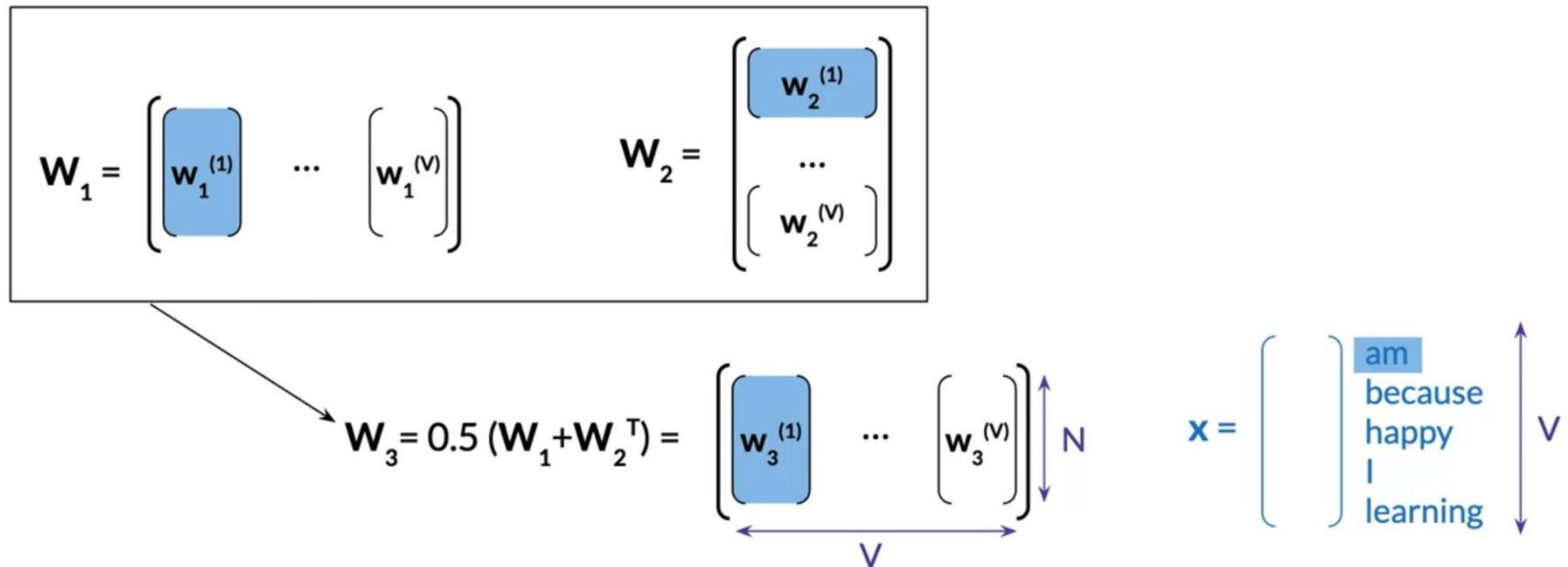
$$\mathbf{W}_2 = \begin{bmatrix} \mathbf{w}^{(1)} \\ \dots \\ \mathbf{w}^{(V)} \end{bmatrix}$$

N

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

V

Extracting word embedding vectors: option 3



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

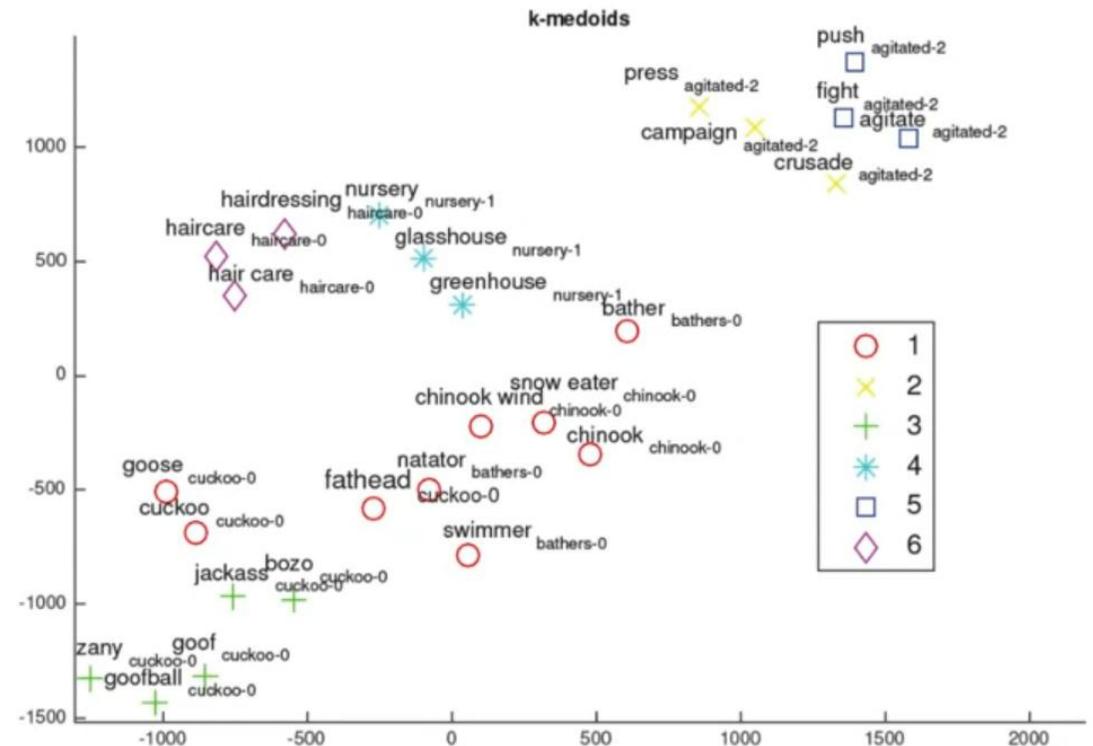
Relationship	Example 1	Example 2	Example 3
France - Paris	Italy: Rome	Japan: Tokyo	Florida: Tallahassee
big - bigger	small: larger	cold: colder	quick: quicker
Miami - Florida	Baltimore: Maryland	Dallas: Texas	Kona: Hawaii
Einstein - scientist	Messi: midfielder	Mozart: violinist	Picasso: painter
Sarkozy - 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)
```