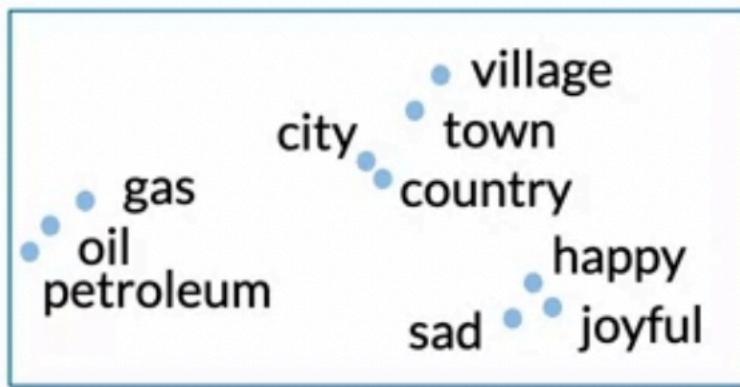


Some basic applications of word embeddings



Semantic analogies
and similarity



Sentiment analysis



Classification of
customer feedback

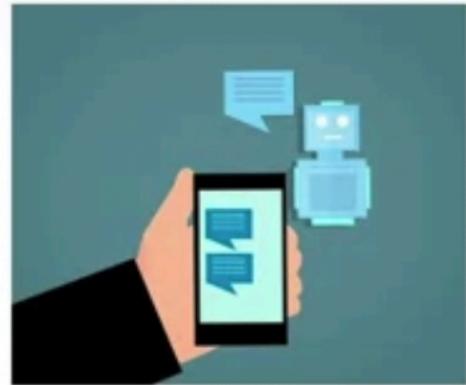
Advanced applications of word embeddings



Machine translation



Information extraction



Question answering

Learning objectives

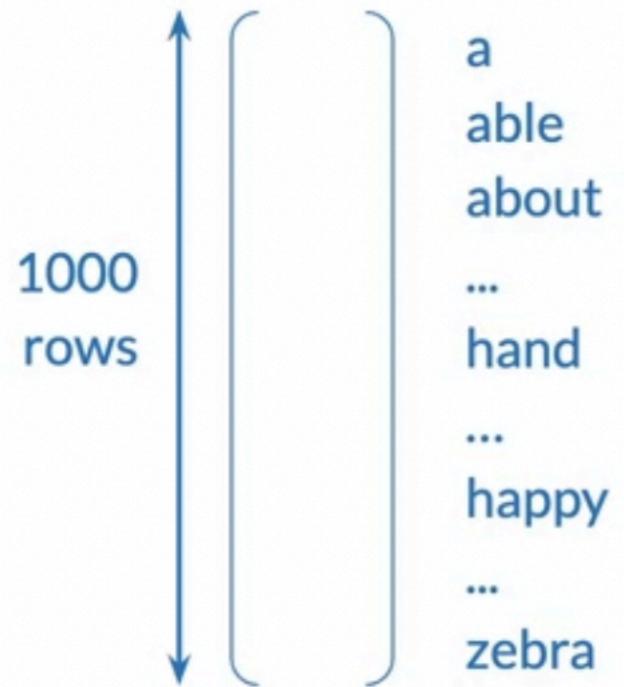
Prerequisite: neural networks

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

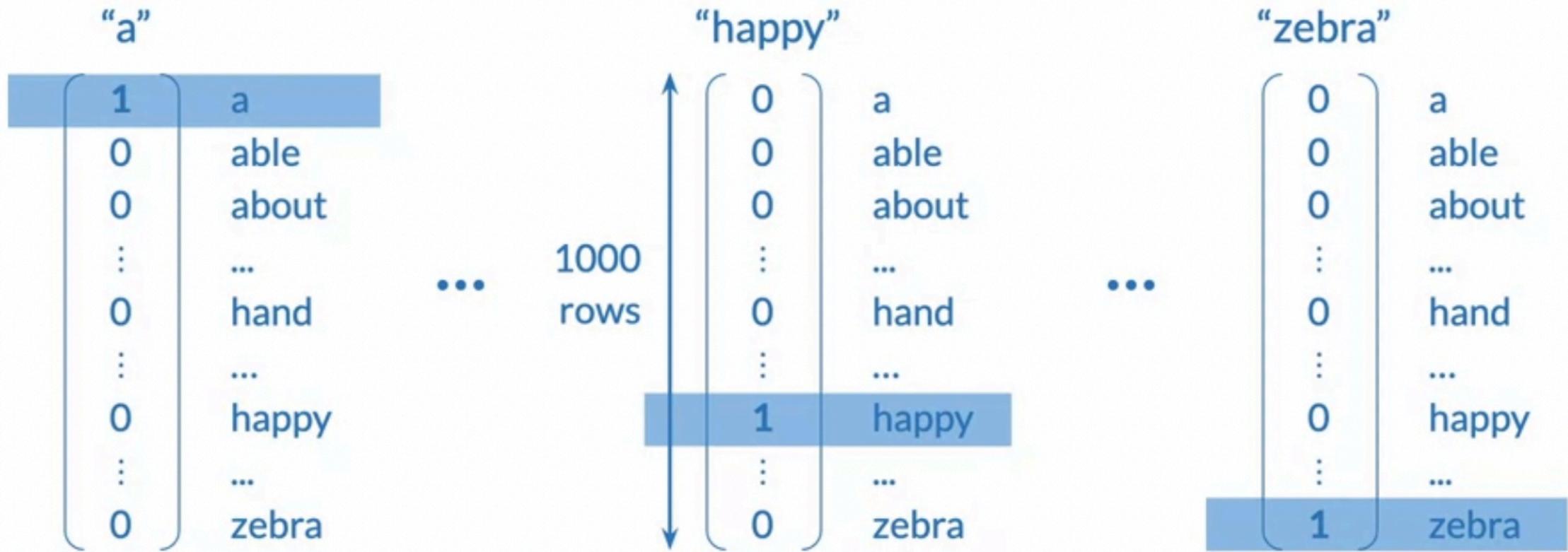
Integers

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

One-hot vectors



One-hot vectors

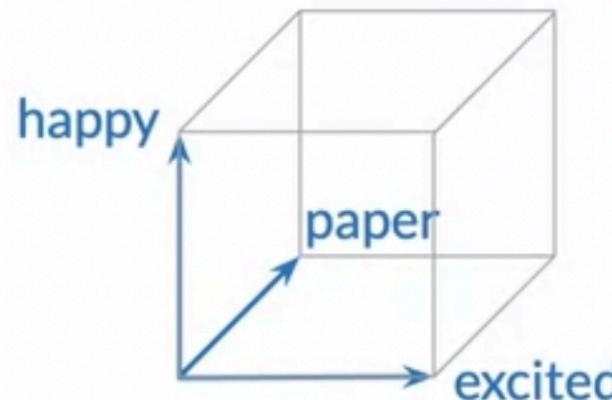
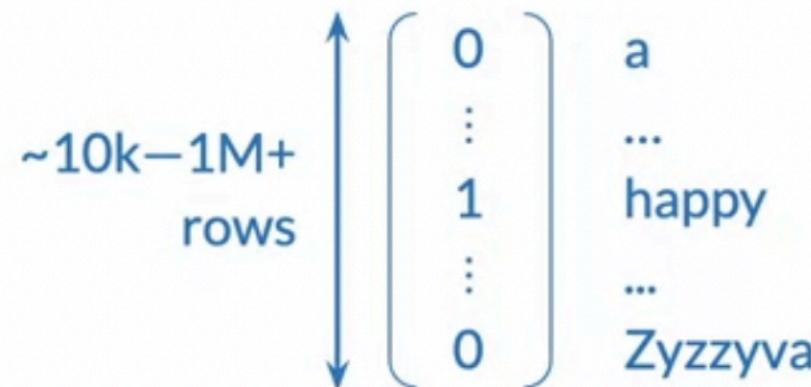


One-hot vectors

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

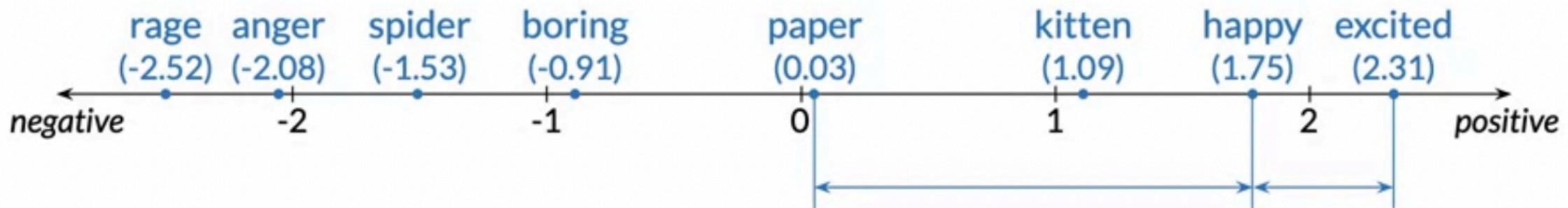
One-hot vectors

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

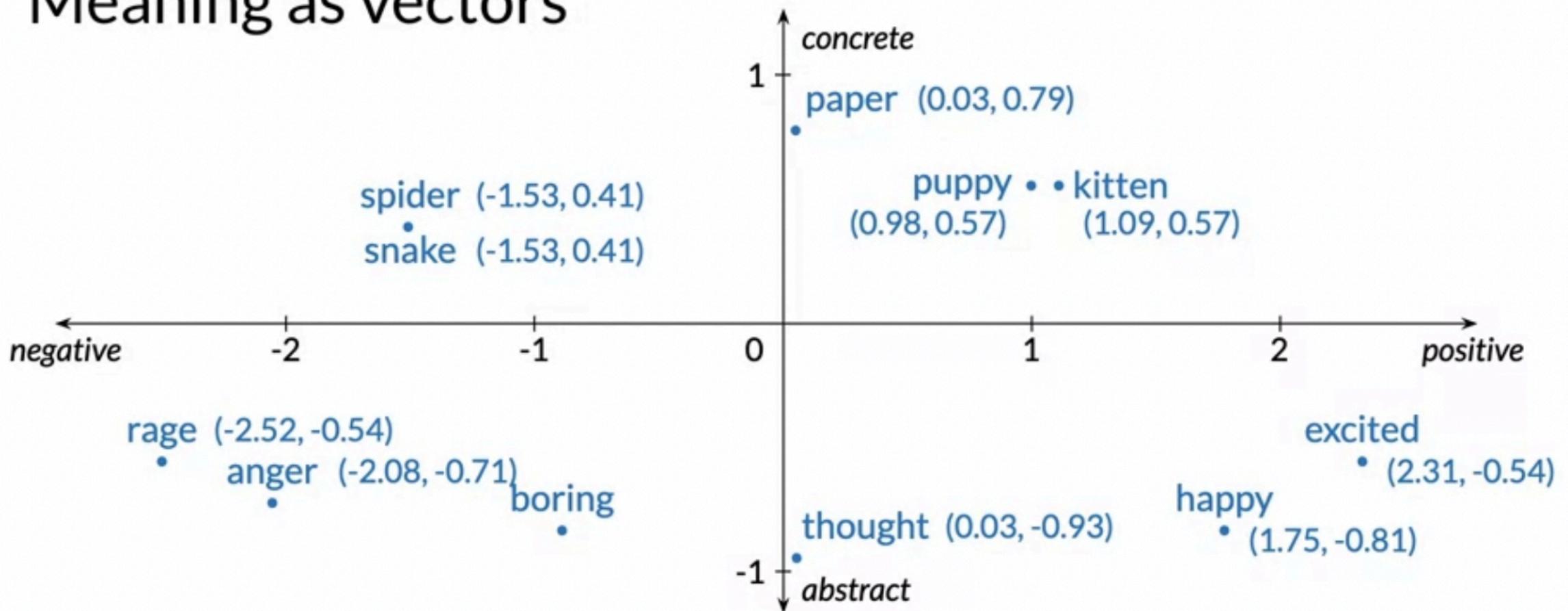


$d(\text{paper}, \text{excited})$
 $= d(\text{paper}, \text{happy})$
 $= d(\text{excited}, \text{happy})$

Meaning as vectors



Meaning as vectors



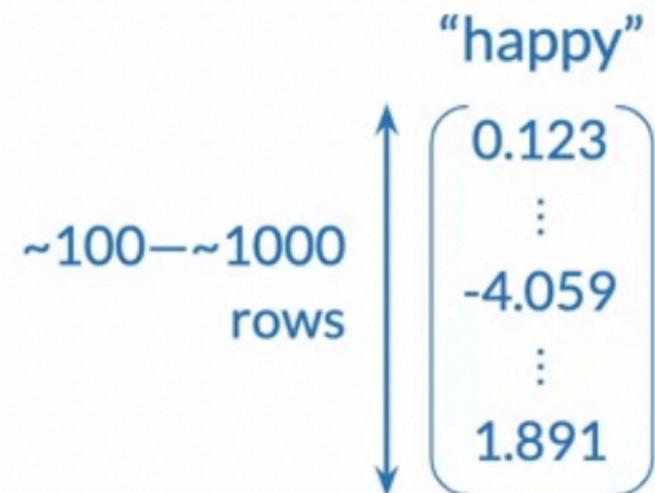
Word embedding vectors

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

$\text{forest} \approx \text{tree}$ $\text{forest} \neq \text{ticket}$

- e.g. analogies

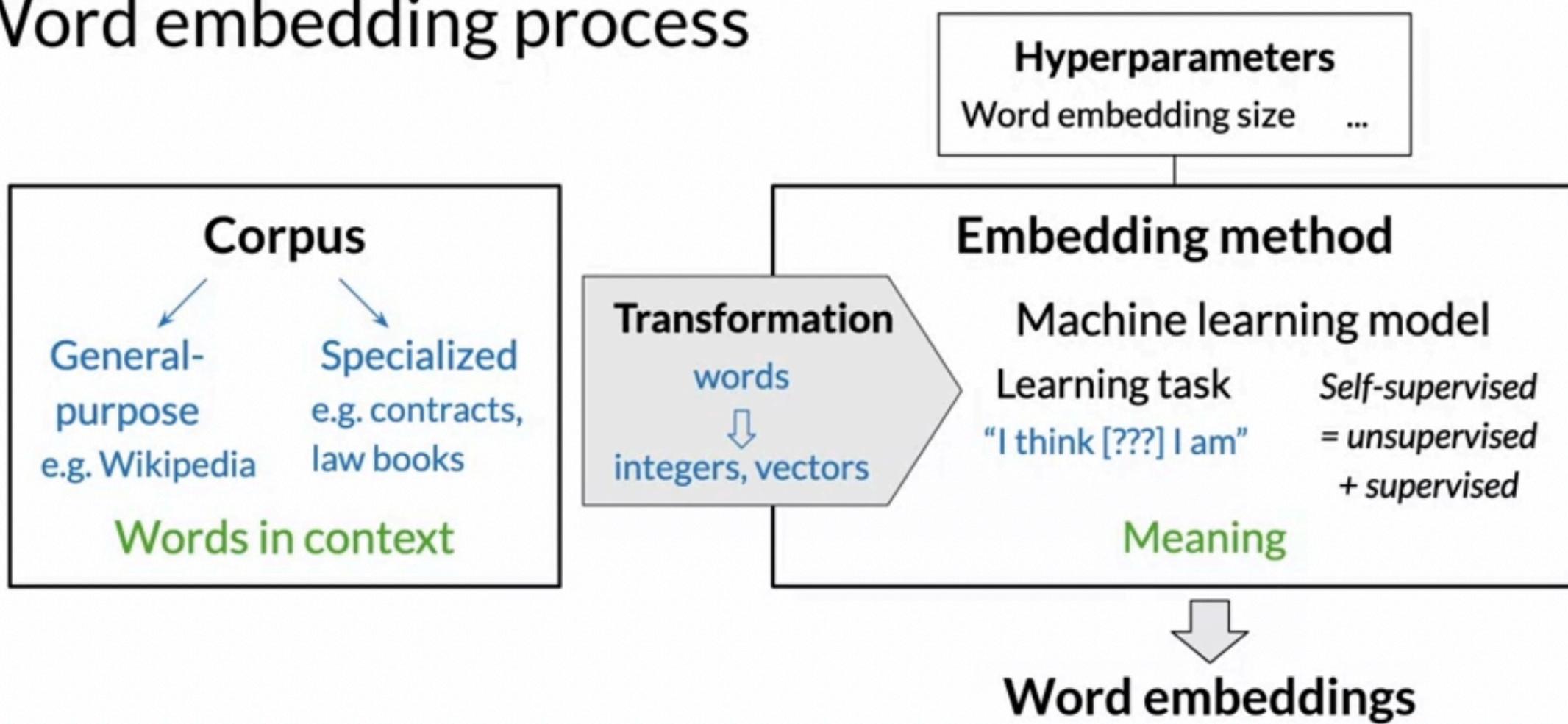
$\text{Paris:France} :: \text{Rome}:?$



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

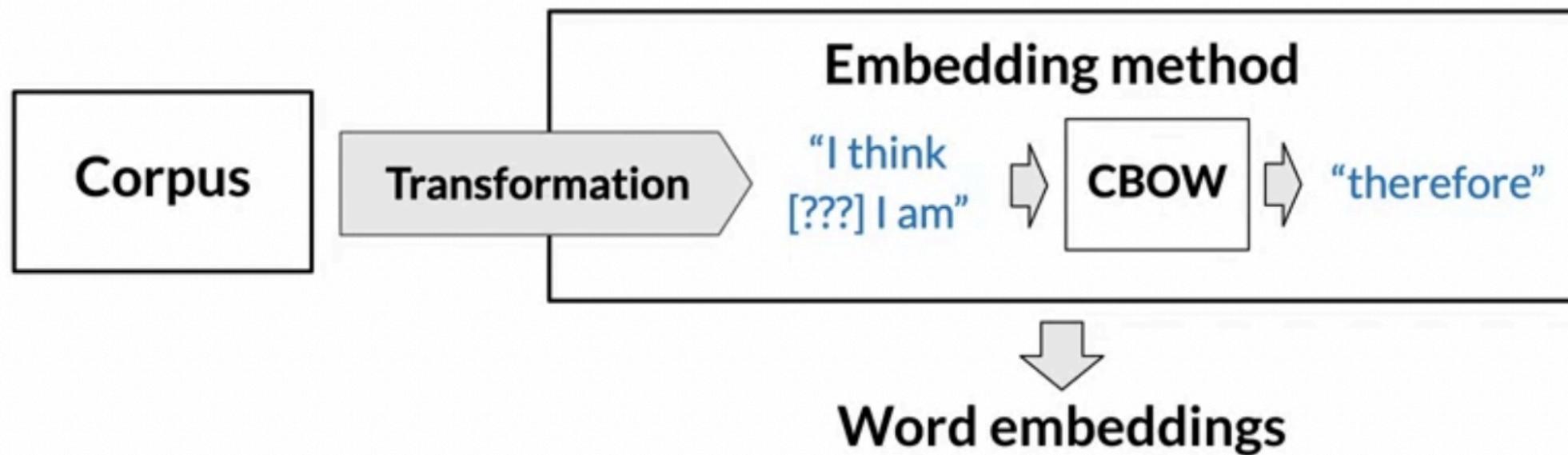
- word2vec (Google, 2013)
 - Continuous bag-of-words (CBOW)
 - Continuous skip-gram / Skip-gram with negative sampling (SGNS)
- Global Vectors (GloVe) (Stanford, 2014)
- fastText (Facebook, 2016)
 - Supports out-of-vocabulary (OOV) words

Advanced word embedding methods

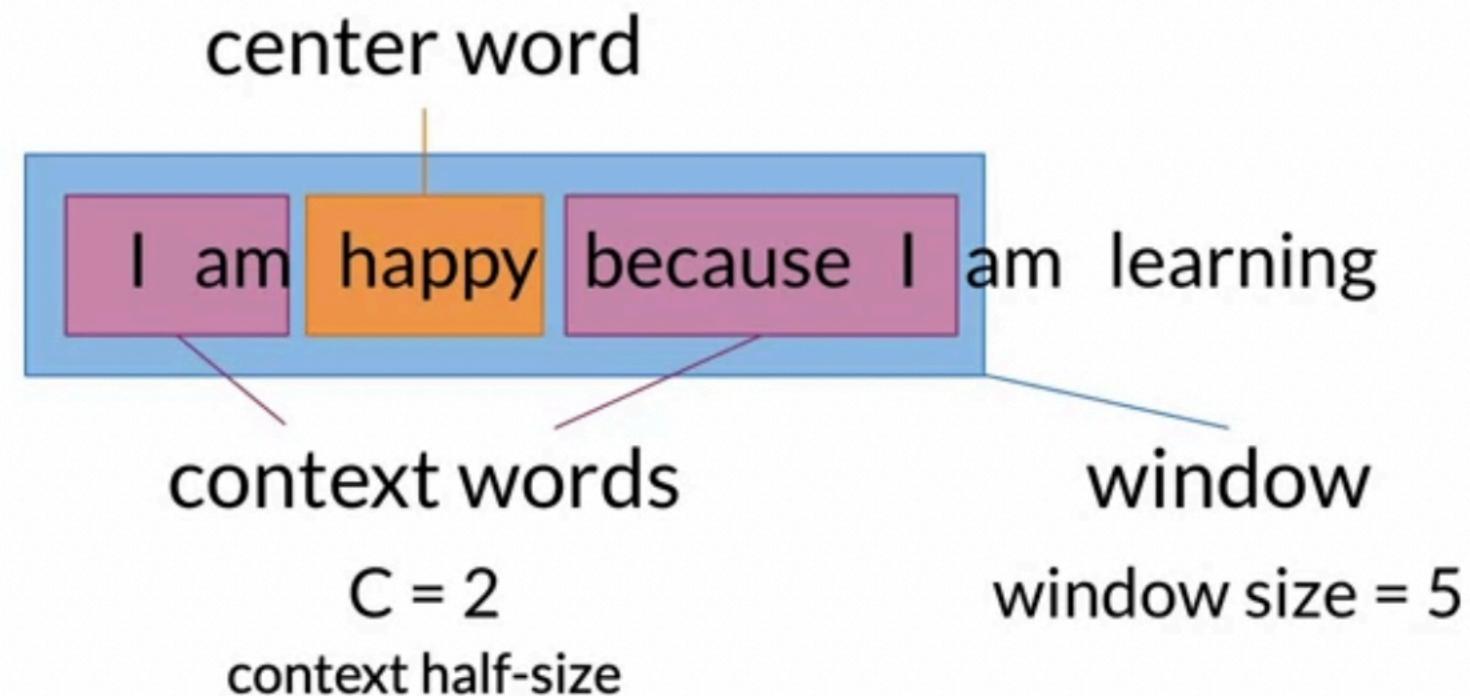
Deep learning, contextual embeddings

- 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



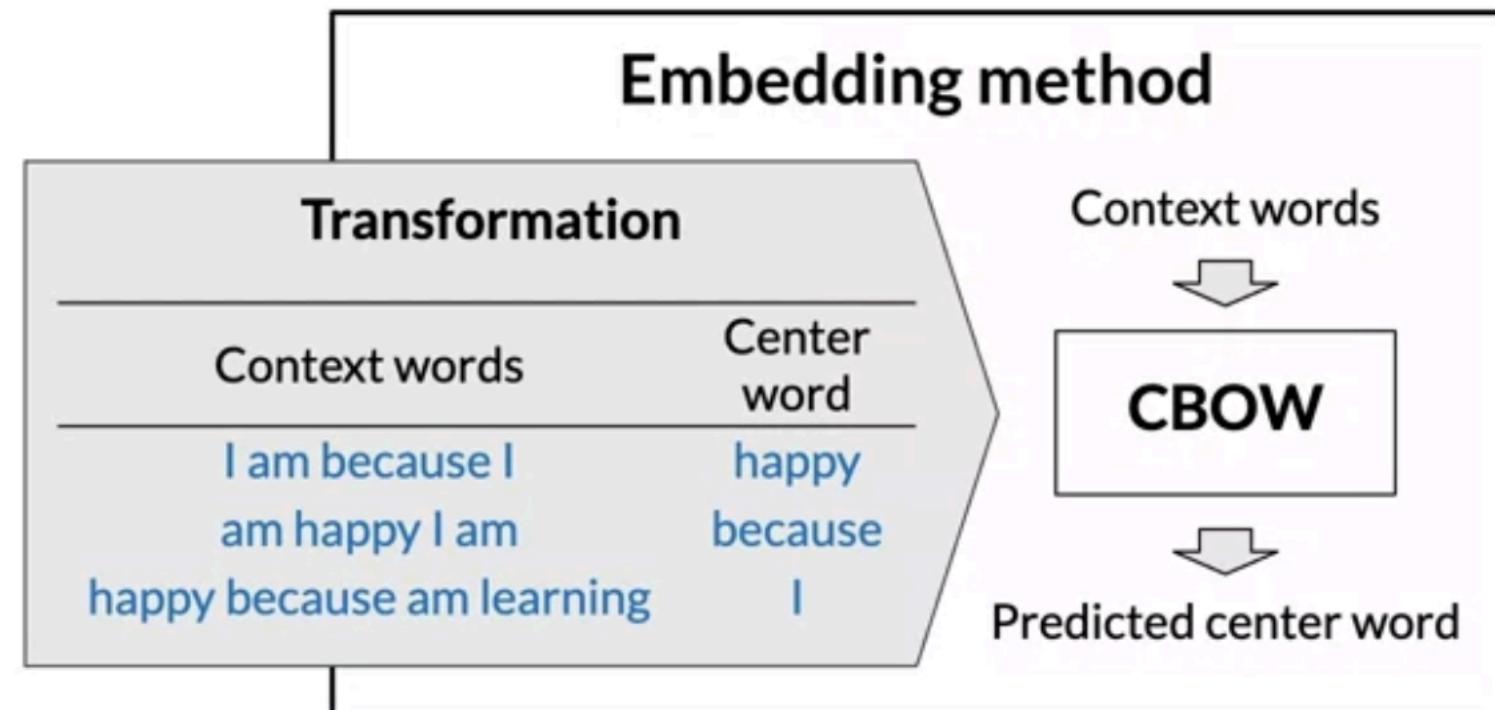
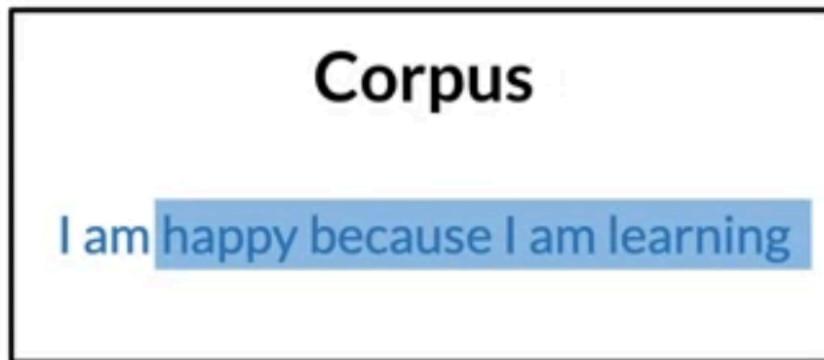
Creating a training example



Corpus

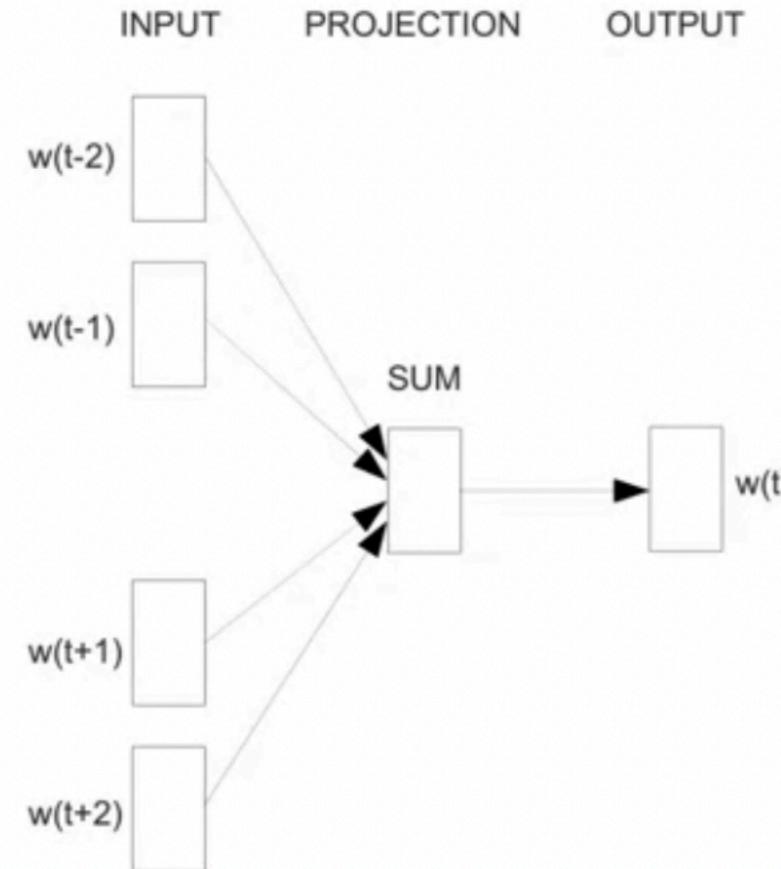
Transformation → CBOW

From corpus to training



Word embeddings

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: 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)
data = nltk.word_tokenize(data) # tokenize string to words

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

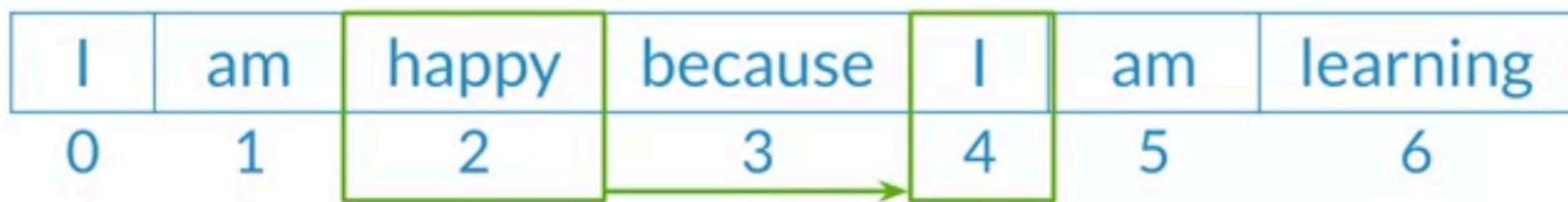
Example in Python: code

```
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
```



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	$\begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$	$\begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{pmatrix}$	$\begin{pmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{pmatrix}$	$\begin{pmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{pmatrix}$	$\begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{pmatrix}$
happy	$\begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{pmatrix}$	$\begin{pmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{pmatrix}$	$\begin{pmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{pmatrix}$	$\begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{pmatrix}$	$\begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{pmatrix}$
I	$\begin{pmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{pmatrix}$	$\begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{pmatrix}$	$\begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{pmatrix}$	$\begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$	$\begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{pmatrix}$
learning	$\begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{pmatrix}$	$\begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{pmatrix}$	$\begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{pmatrix}$	$\begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{pmatrix}$	$\begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{pmatrix}$

Transforming context words into vectors

Average of individual one-hot vectors

$$\left[\begin{array}{c} I \\ am \\ because \\ happy \\ I \\ learning \end{array} \right] + \left[\begin{array}{c} 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{array} \right] + \left[\begin{array}{c} 1 \\ 0 \\ 1 \\ 0 \\ 0 \end{array} \right] + \left[\begin{array}{c} 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{array} \right] / 4 = \left[\begin{array}{c} 0.25 \\ 0.25 \\ 0 \\ 0.5 \\ 0 \end{array} \right]$$

I am because I

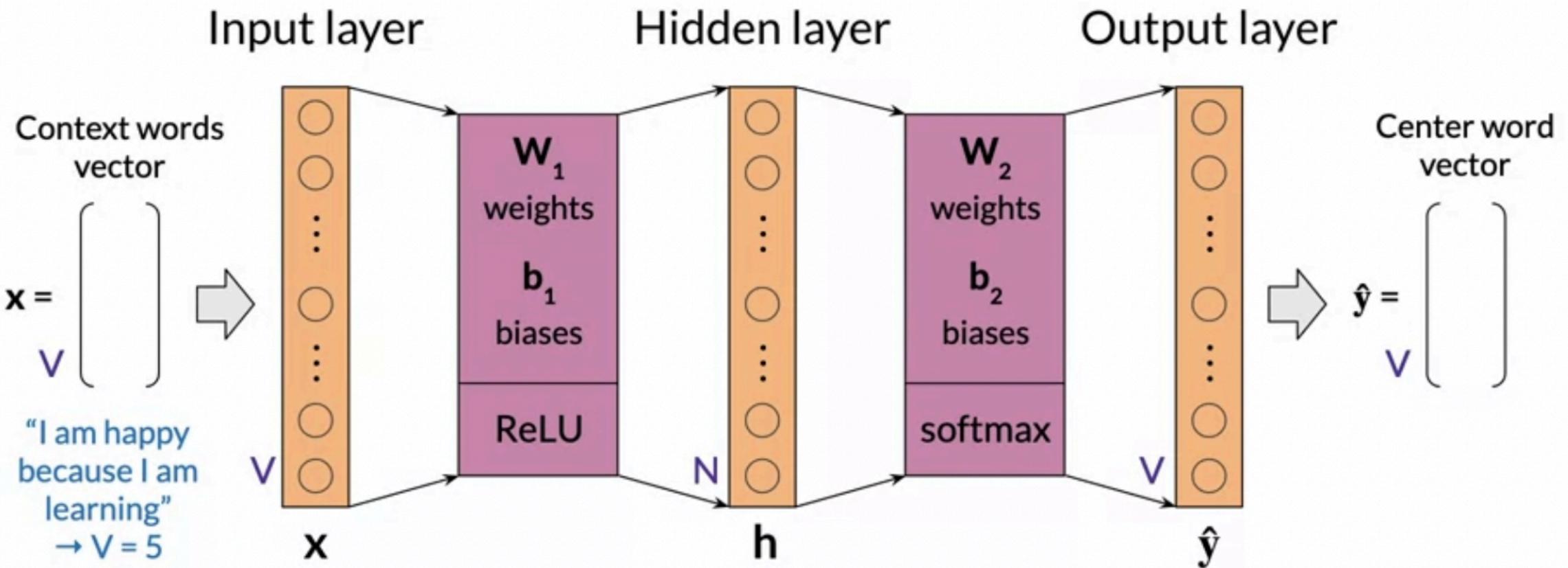
Final prepared training set

Context words	Context words vector	Center word	Center word vector
<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]

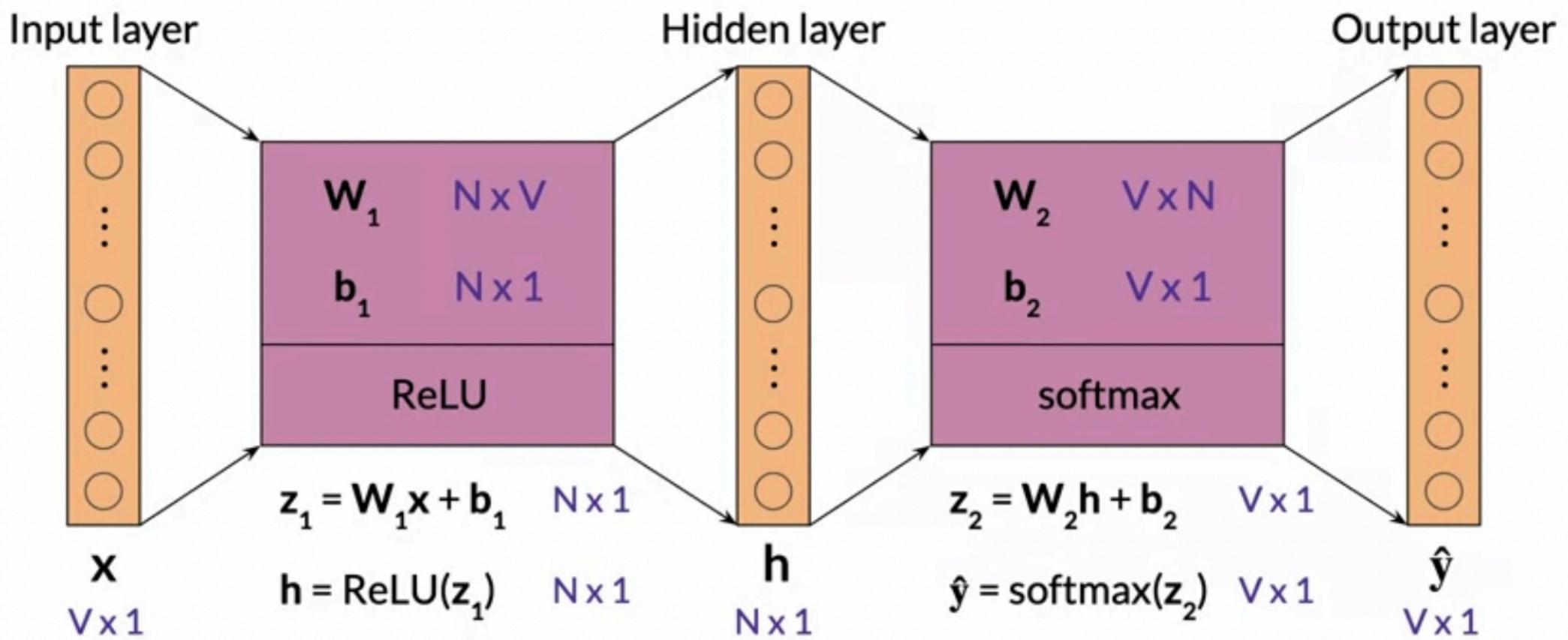
Architecture of the CBOW model

Hyperparameters

N: Word embedding size ...



Dimensions (single input)



Dimensions (single input)

Column vectors

$$\mathbf{z}_1 = \mathbf{W}_1 \mathbf{x} + \mathbf{b}_1$$

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

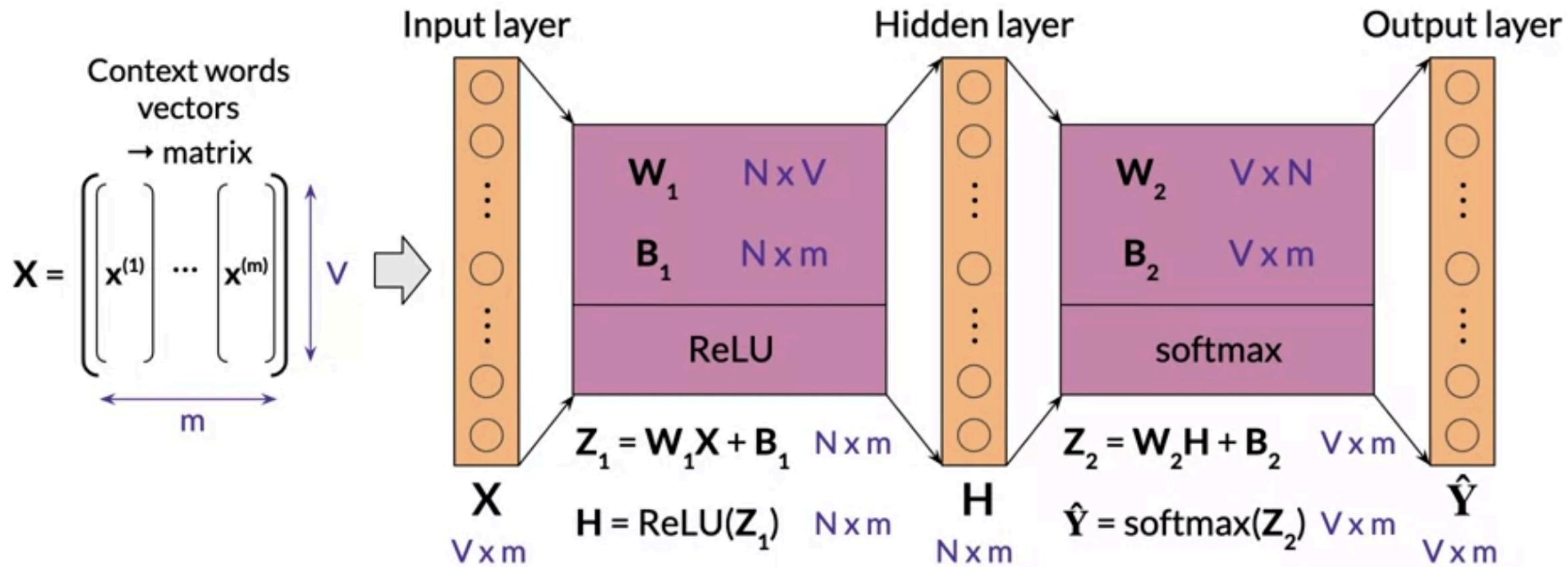
Row vectors

$$\mathbf{z}_1 = \mathbf{x} \mathbf{W}_1^T + \mathbf{b}_1$$

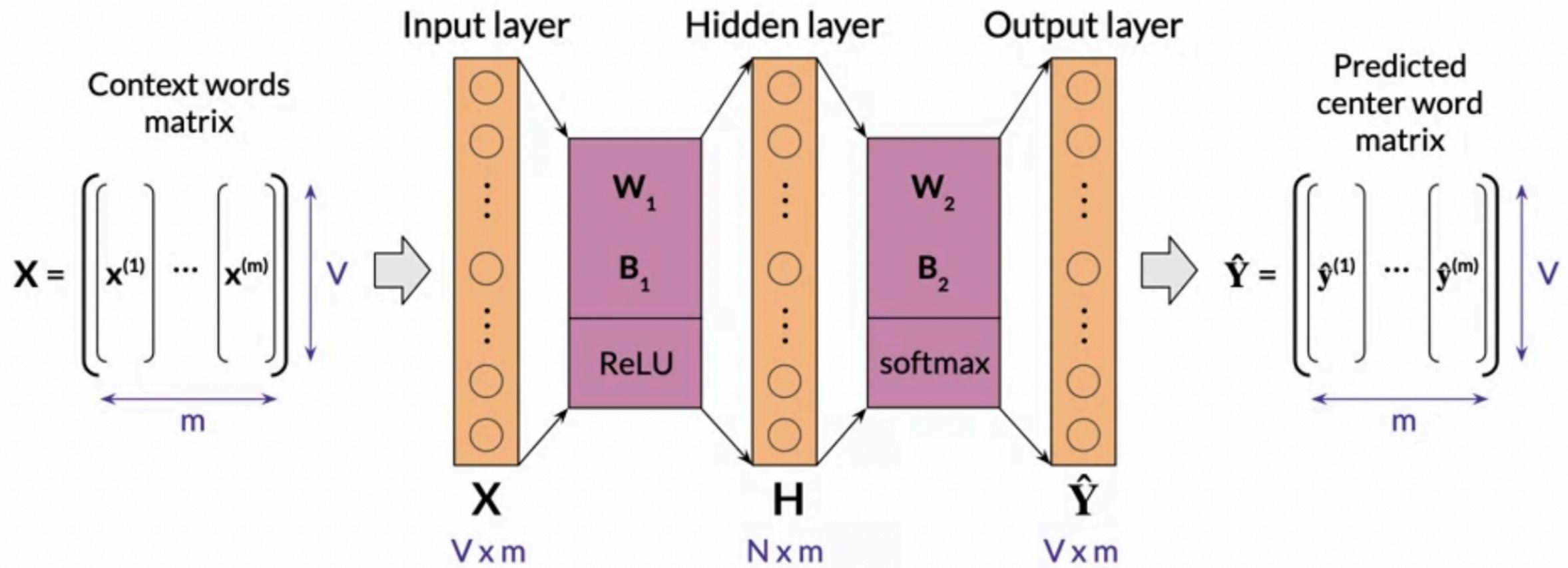
$$\mathbf{b}_1 = \begin{bmatrix} \quad \\ 1 \times N \end{bmatrix} \quad \mathbf{W}_1 = \begin{bmatrix} \quad \\ \quad \end{bmatrix}_{N \times V} \quad \mathbf{b}_1 = \begin{bmatrix} \quad \\ 1 \times N \end{bmatrix}$$
$$\mathbf{x} = \begin{bmatrix} \quad \\ 1 \times V \end{bmatrix}$$

Dimensions (batch input)

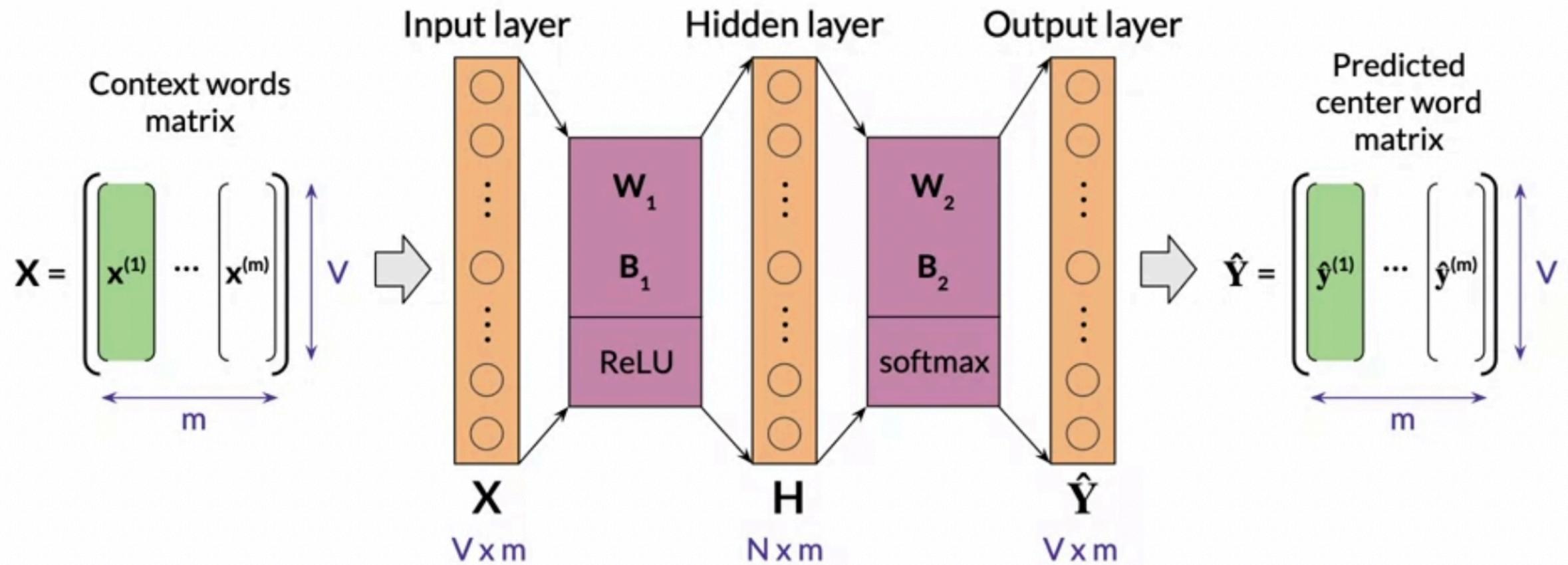
$$\begin{bmatrix} \mathbf{b}_1 \end{bmatrix} \rightarrow \mathbf{B}_1 = \begin{bmatrix} \begin{bmatrix} \mathbf{b}_1 \end{bmatrix} & \dots & \begin{bmatrix} \mathbf{b}_1 \end{bmatrix} \end{bmatrix} \underset{m}{\underbrace{\qquad\qquad\qquad}} \underset{N}{\updownarrow} \text{broadcasting}$$



Dimensions (batch input)

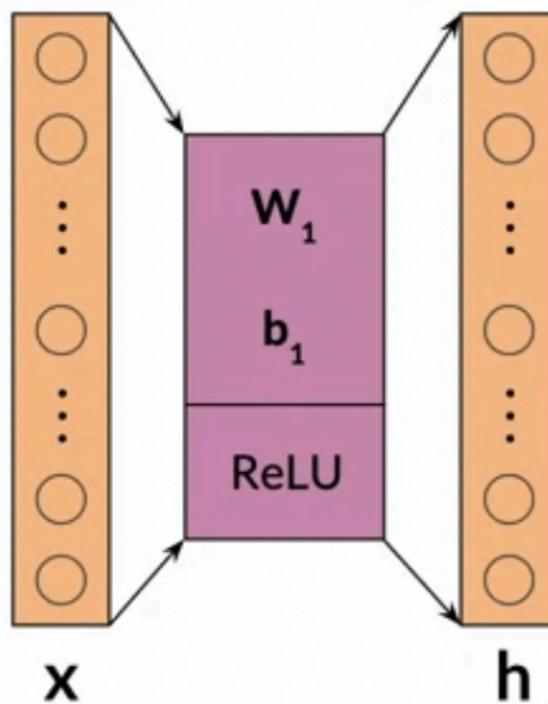


Dimensions (batch input)



Rectified Linear Unit (ReLU)

Input layer Hidden layer

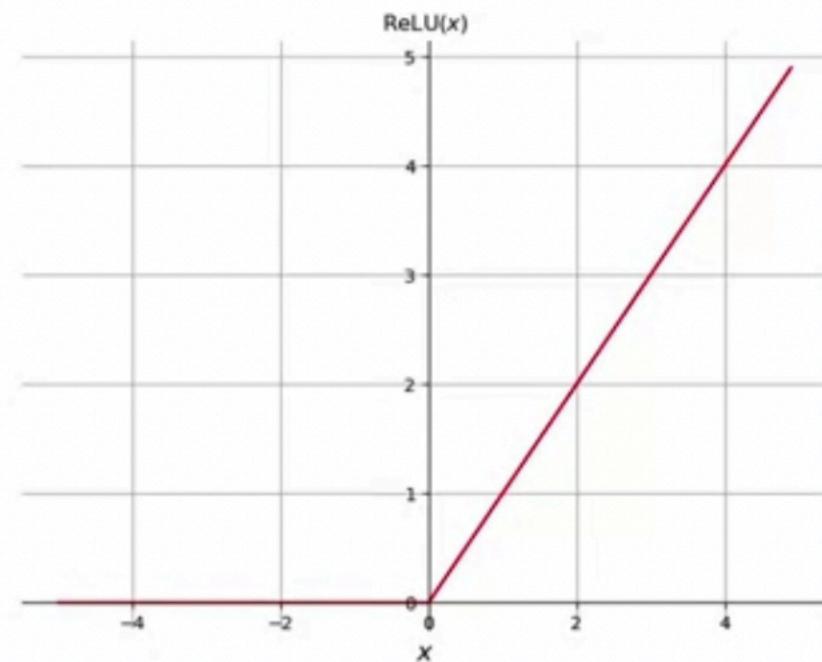


$$z_1 = \mathbf{W}_1 \mathbf{x} + \mathbf{b}_1$$

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

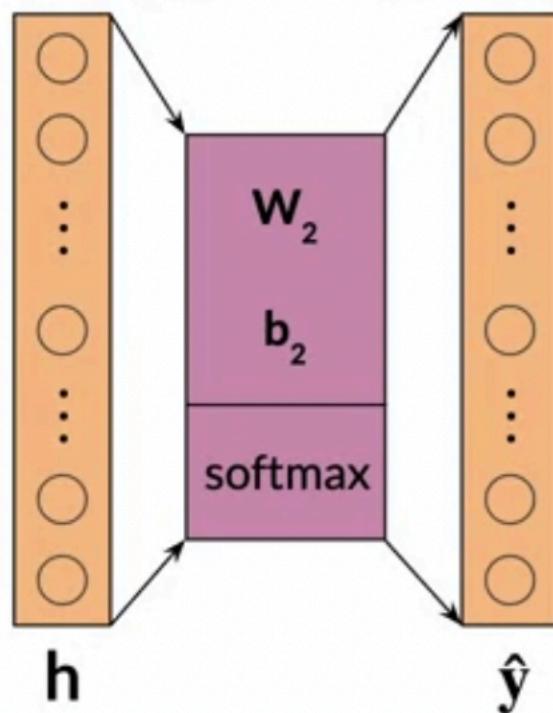
A diagram showing the calculation of z_1 and \mathbf{h} . On the left, a blue vertical vector labeled z_1 is shown with values 5.1 , -0.3 , \vdots , -4.6 , and 0.2 . An arrow labeled "ReLU" points to the right, where a blue vertical vector labeled \mathbf{h} is shown with values 0 , \vdots , and 0 .

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



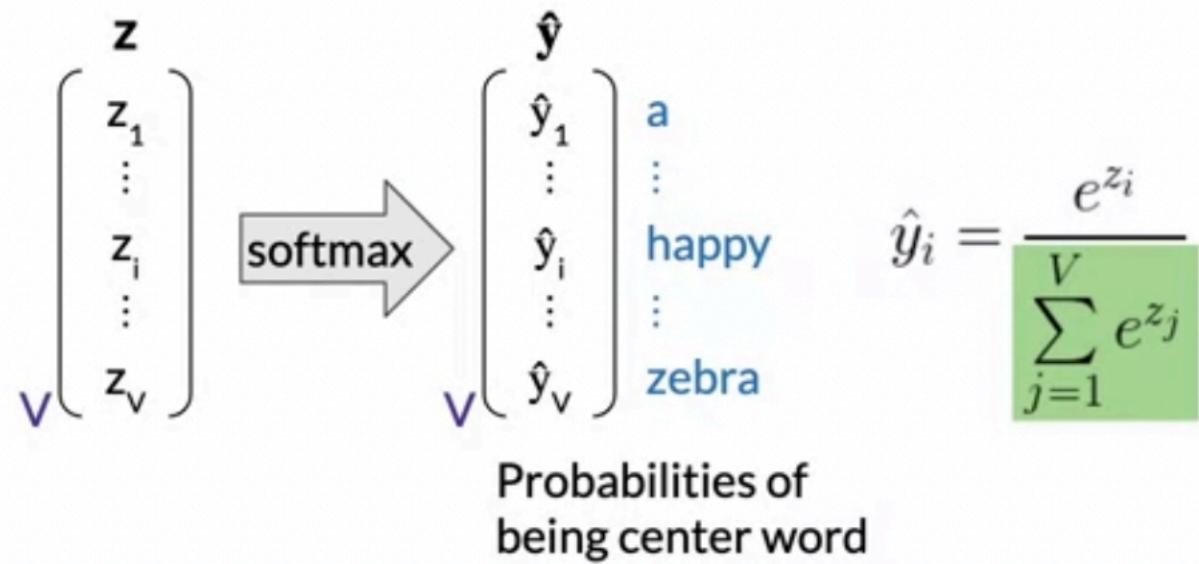
Softmax

Hidden layer Output layer



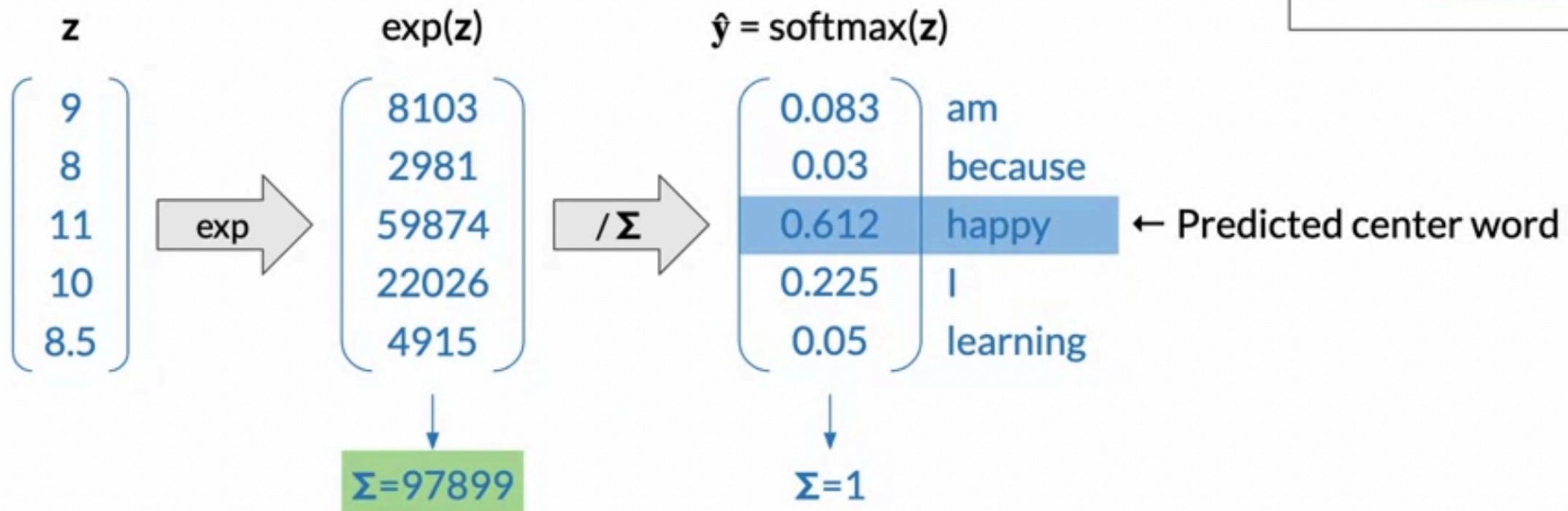
$$z = \mathbf{W}_2 \mathbf{h} + \mathbf{b}_2$$

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

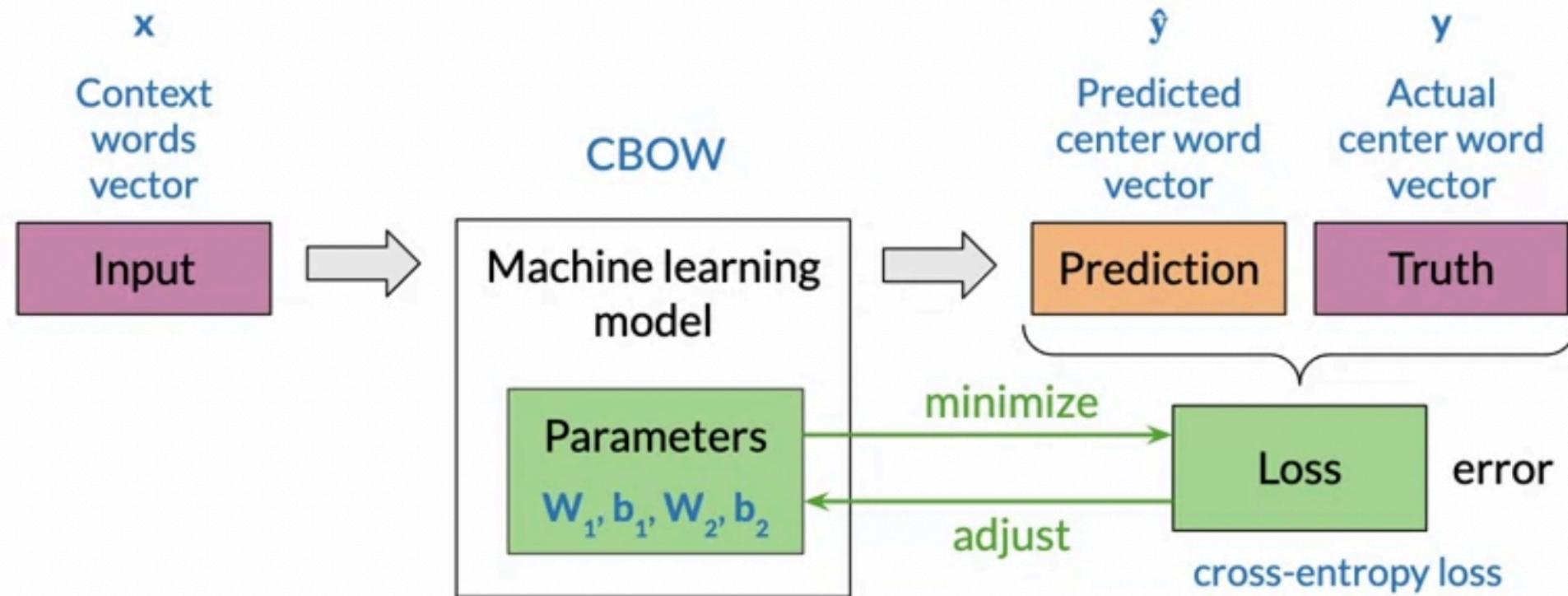


Softmax: example

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



LOSS

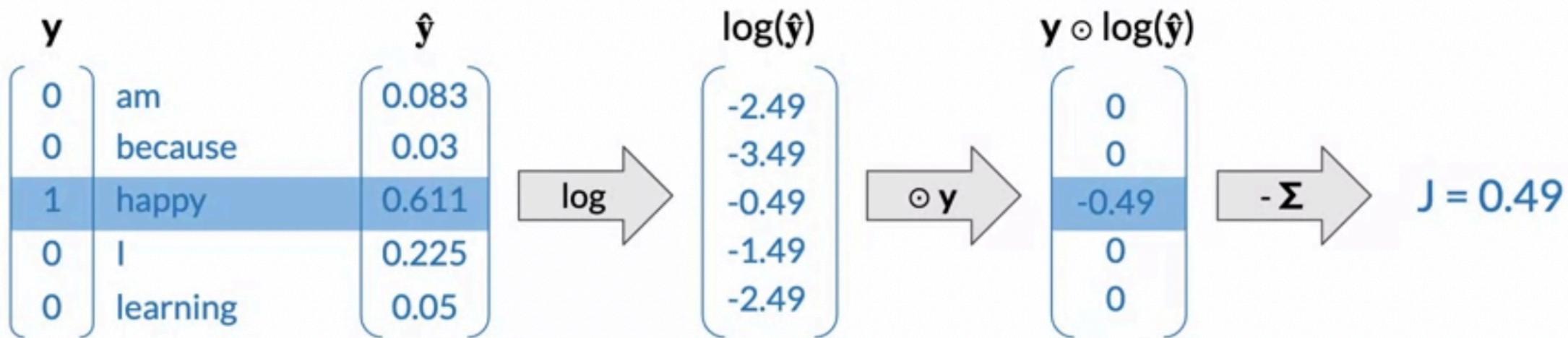


Cross-entropy loss

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

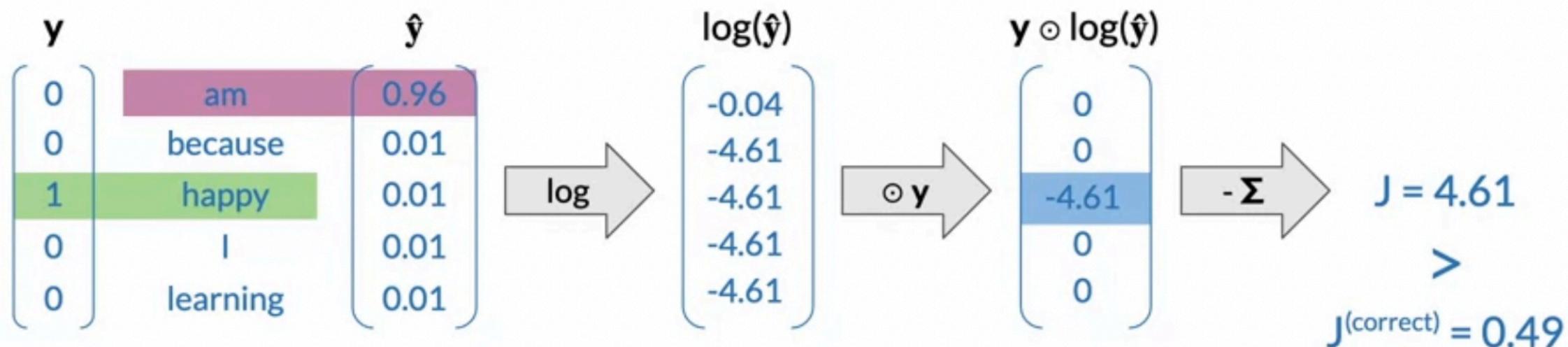
Actual $\mathbf{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_v \end{pmatrix}$	Predicted $\hat{\mathbf{y}} = \begin{pmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_v \end{pmatrix}$
---	--

I am happy because I am learning



Cross-entropy loss

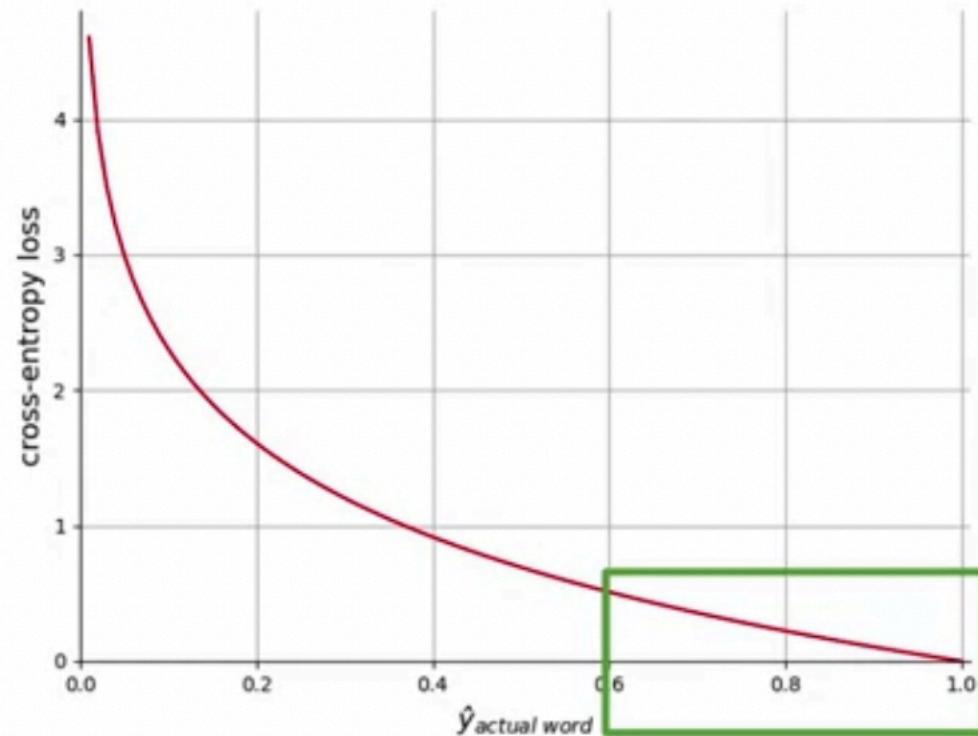
$$J = - \sum_{k=1}^V y_k \log \hat{y}_k$$



Cross-entropy loss

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

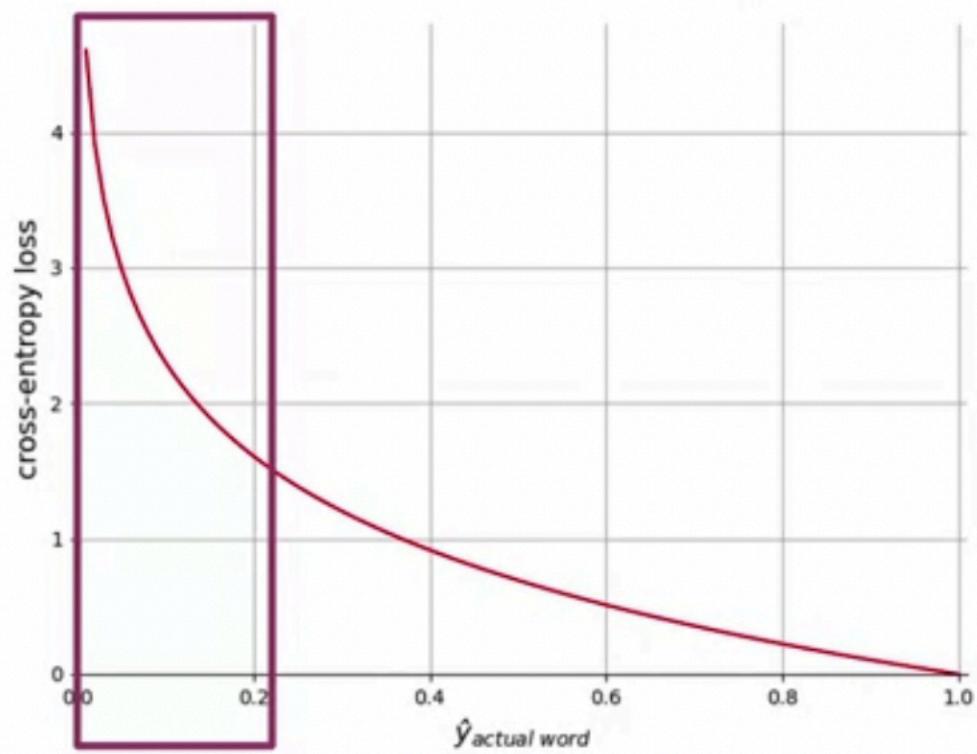
$$\begin{array}{c} \mathbf{y} \\ \left[\begin{array}{c} 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{array} \right] \end{array} \quad \begin{array}{c} \text{am} \\ \text{because} \\ \text{happy} \\ \text{I} \\ \text{learning} \end{array} \quad \begin{array}{c} \hat{\mathbf{y}} \\ \left[\begin{array}{c} 0.96 \\ 0.01 \\ 0.01 \\ 0.01 \\ 0.01 \end{array} \right] \end{array} \rightarrow J = 4.61$$



Cross-entropy loss

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

$$\begin{array}{c} \mathbf{y} \\ \left[\begin{array}{c} 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{array} \right] \end{array} \quad \begin{array}{c} \text{am} \\ \text{because} \\ \text{happy} \\ \text{I} \\ \text{learning} \end{array} \quad \begin{array}{c} \hat{\mathbf{y}} \\ \left[\begin{array}{c} 0.96 \\ 0.01 \\ 0.01 \\ 0.01 \\ 0.01 \end{array} \right] \end{array} \rightarrow J = 4.61$$



Training process

- Forward propagation
- Cost
- Backpropagation and gradient descent

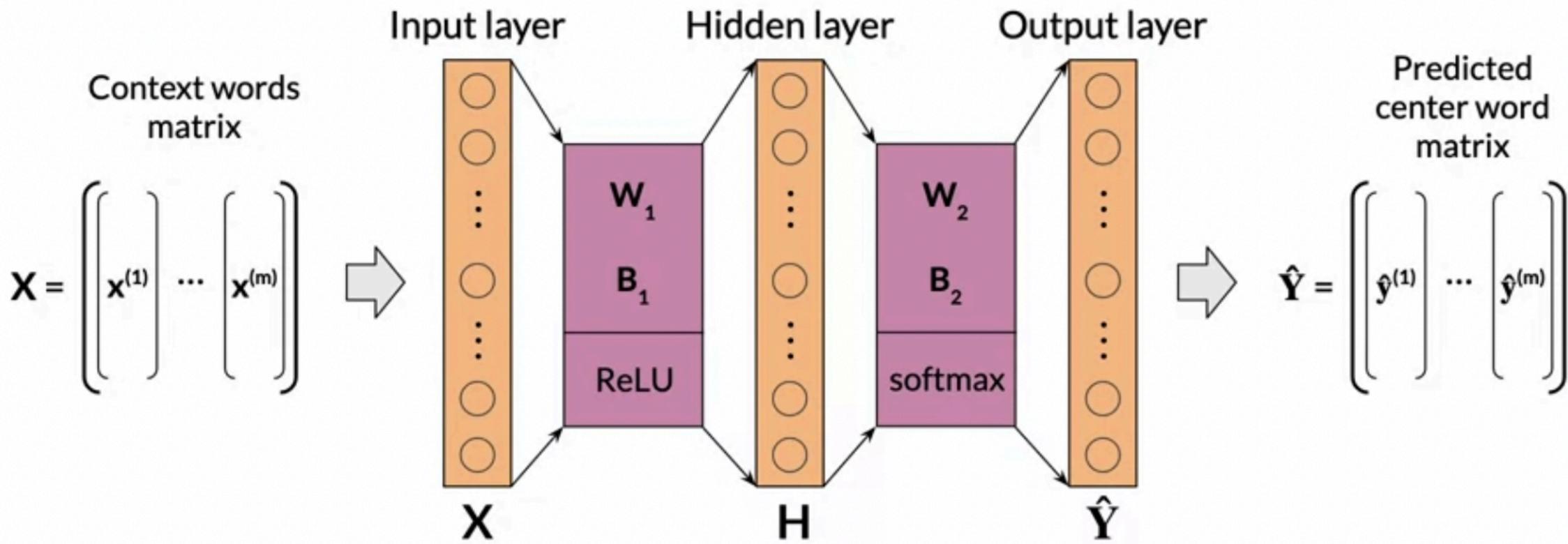
Forward propagation

$$Z_1 = W_1 X + B_1$$

$$H = \text{ReLU}(Z_1)$$

$$Z_2 = W_2 H + B_2$$

$$\hat{Y} = \text{softmax}(Z_2)$$



Cost

Cost: mean of losses

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

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

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

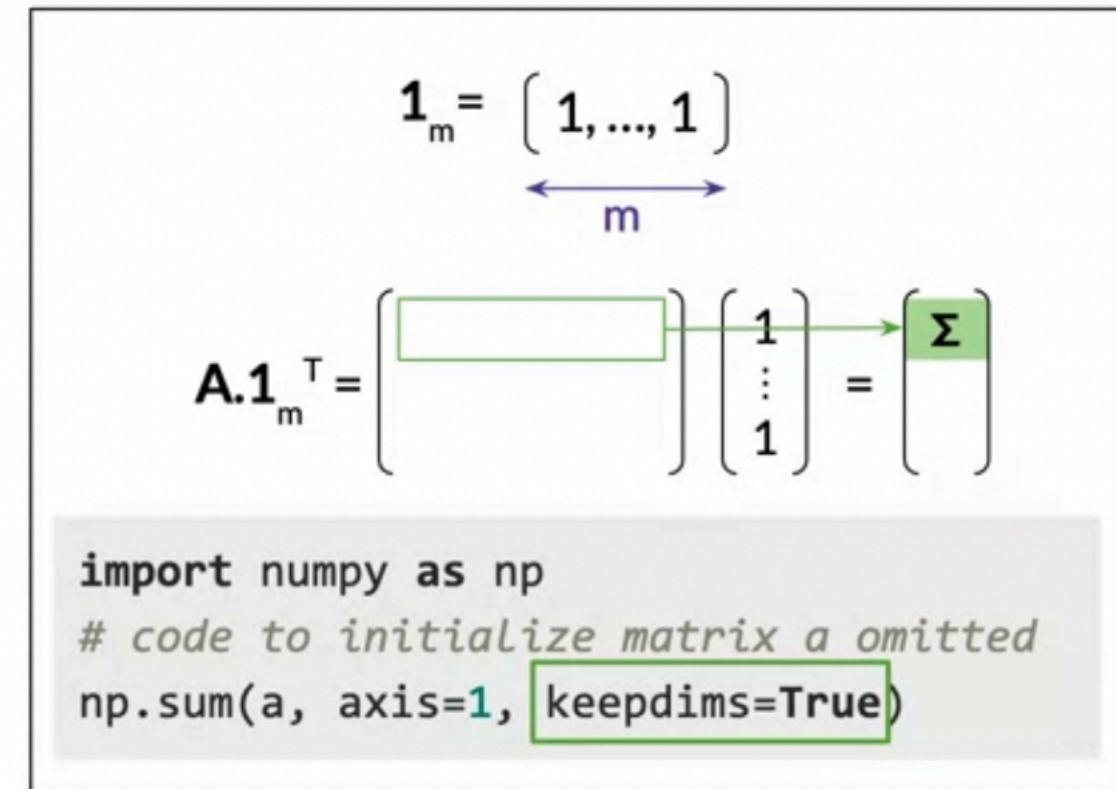
Backpropagation

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

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

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

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



Gradient descent

Hyperparameter: learning rate α

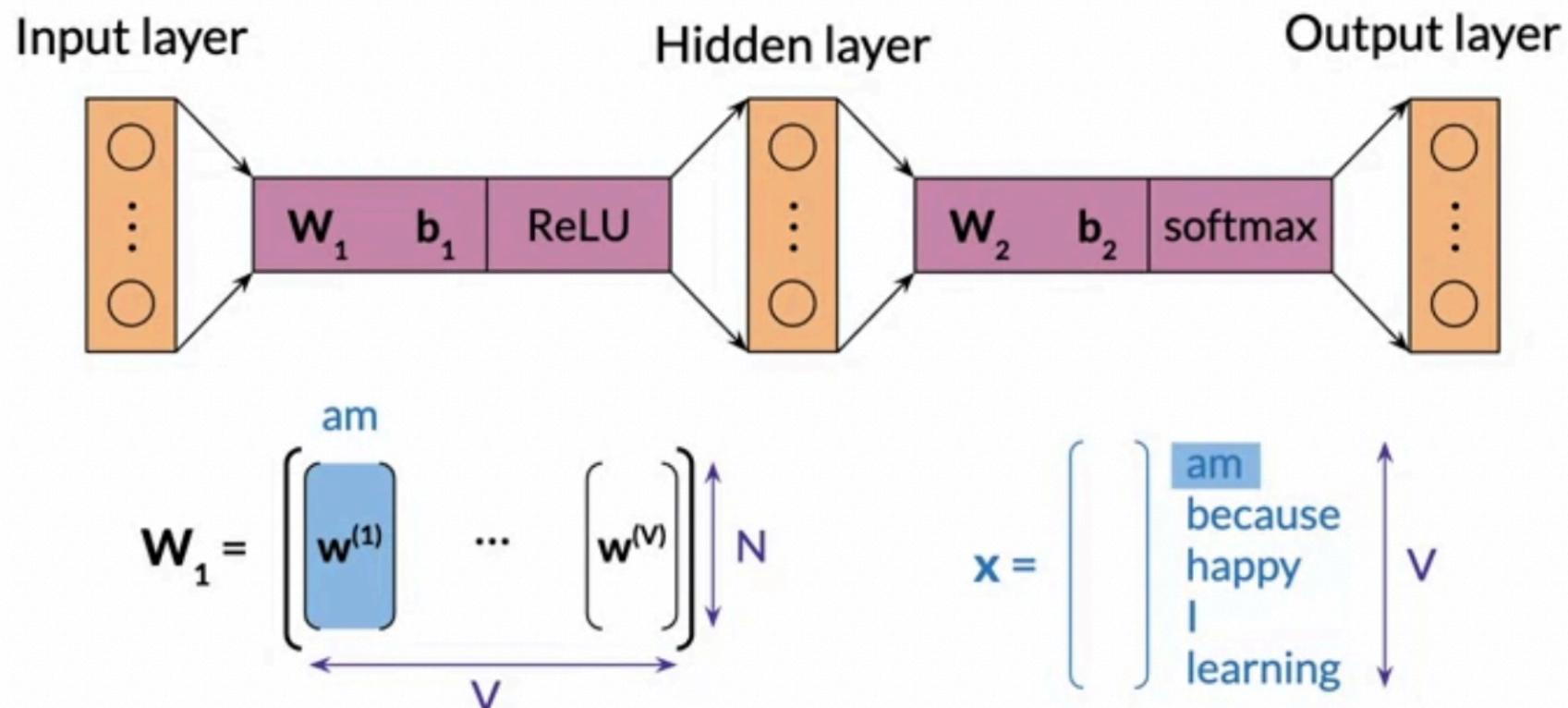
$$\mathbf{W}_1 := \mathbf{W}_1 - \alpha \frac{\partial J_{batch}}{\partial \mathbf{W}_1}$$

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

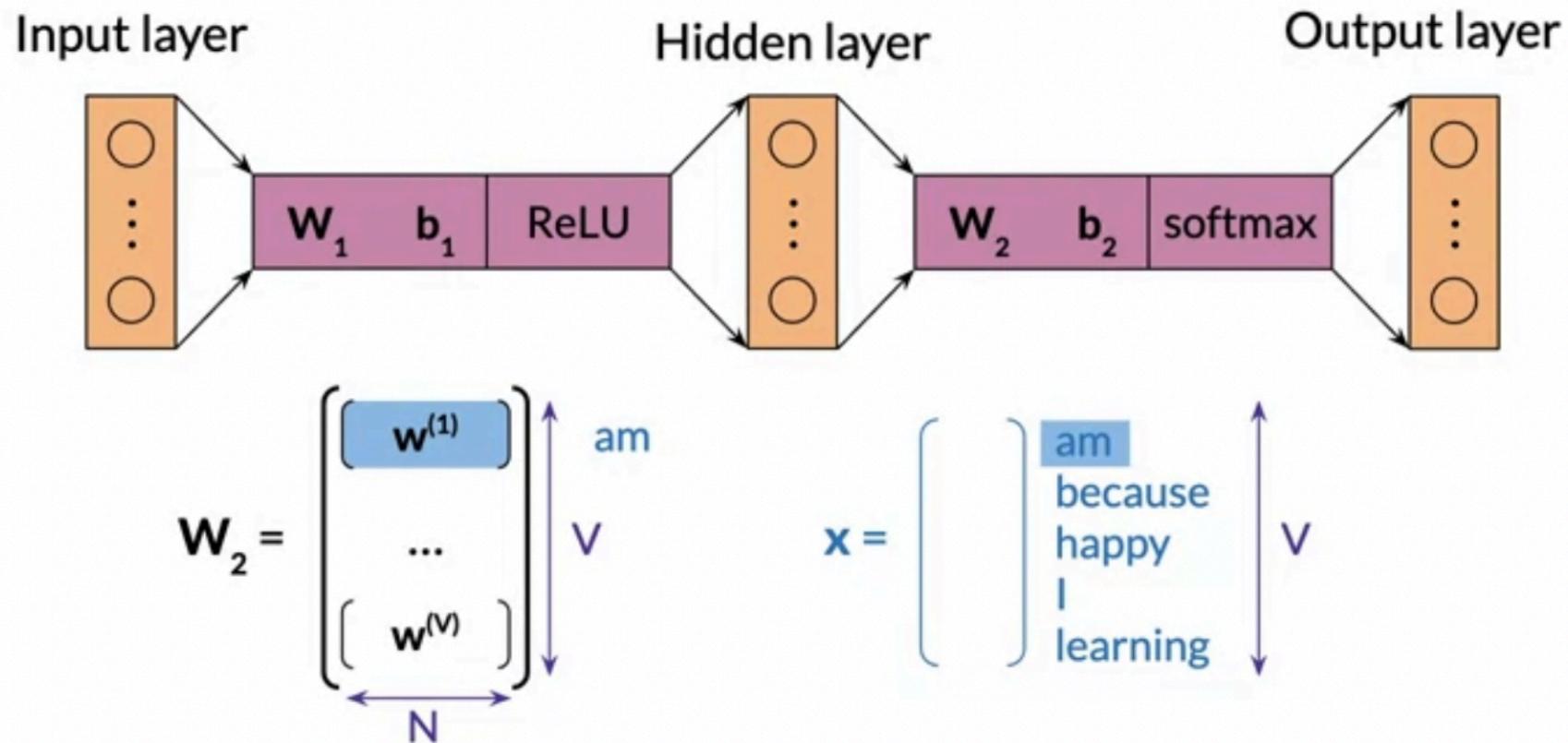
$$\mathbf{b}_1 := \mathbf{b}_1 - \alpha \frac{\partial J_{batch}}{\partial \mathbf{b}_1}$$

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

Extracting word embedding vectors: option 1



Extracting word embedding vectors: option 2



Extracting word embedding vectors: option 3

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

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

$$\mathbf{W}_3 = 0.5 (\mathbf{W}_1 + \mathbf{W}_2^T) = \begin{bmatrix} \mathbf{w}_3^{(1)} & \dots & \mathbf{w}_3^{(V)} \end{bmatrix}$$

$\xleftarrow{V} \quad \xrightarrow{N}$

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

$\uparrow V \quad \downarrow V$

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

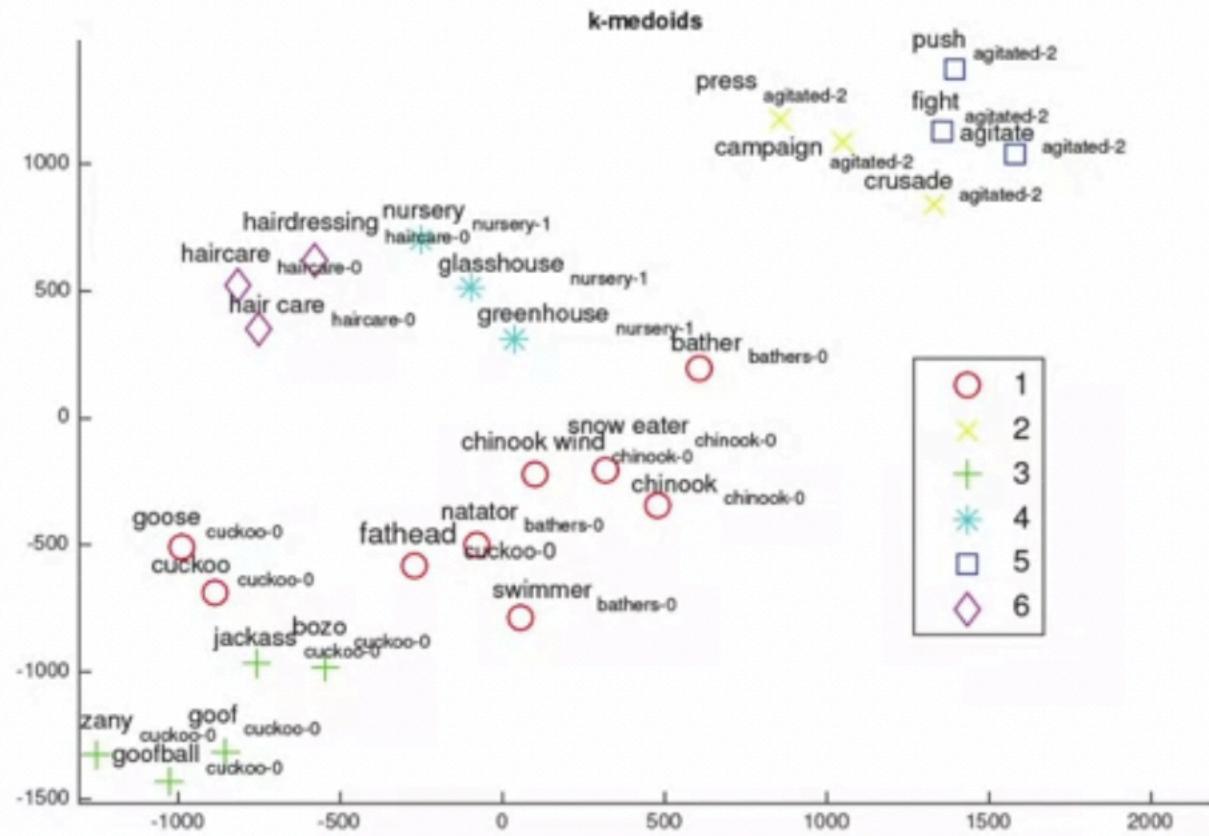
Relationship	Example 1	Example 2	Example 3
France - Paris	Italy: Rome	Japan: Tokyo	Florida: Tallahassee
big - bigger	small: larger	cold: colder	quick: quicker
Miami - Florida	Baltimore: Maryland	Dallas: Texas	Kona: Hawaii
Einstein - scientist	Messi: midfielder	Mozart: violinist	Picasso: painter
Sarkozy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

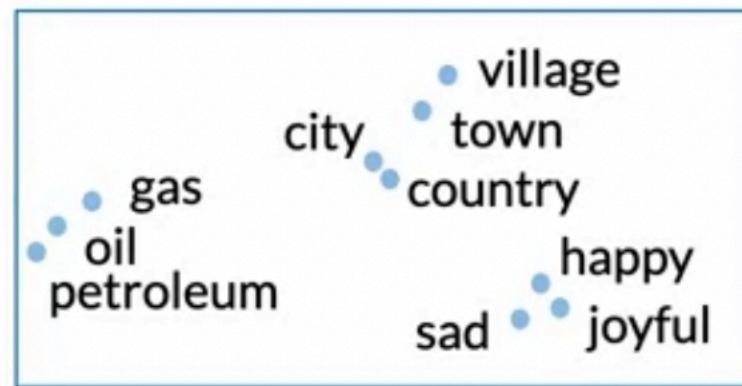
Source: Michael Zhai, Johnny Tan, and Jinho D. Choi. 2016. [Intrinsic and extrinsic evaluations of word embeddings](#)



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

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)*