

Some basic applications of word embeddings



Semantic analogies
and similarity



Sentiment analysis



Classification of
customer feedback

Advanced applications of word embeddings



Machine translation

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

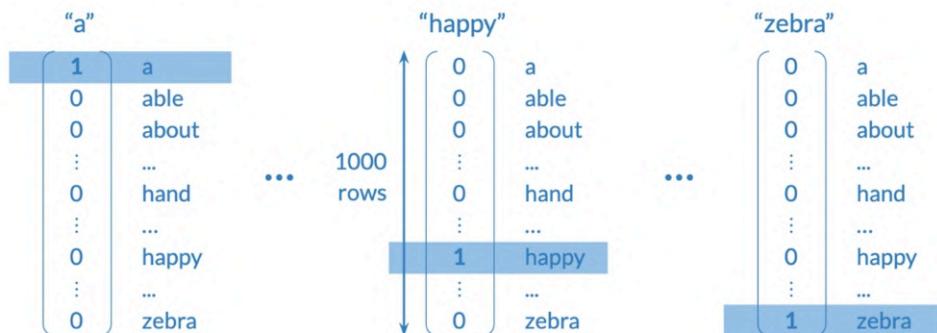
+ Simple

- Ordering: little semantic sense

hand 615 < happy 621 < zebra 1000

?! ?!

One-hot vectors



One-hot vectors

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

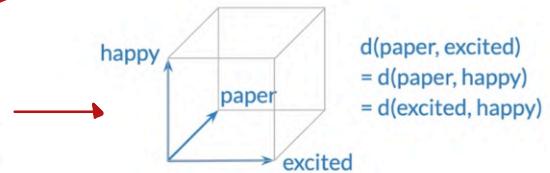
One-hot vectors

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

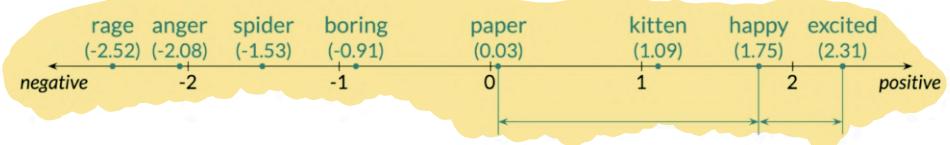
$\sim 10k - 1M+$ rows

$$\begin{pmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{pmatrix}$$

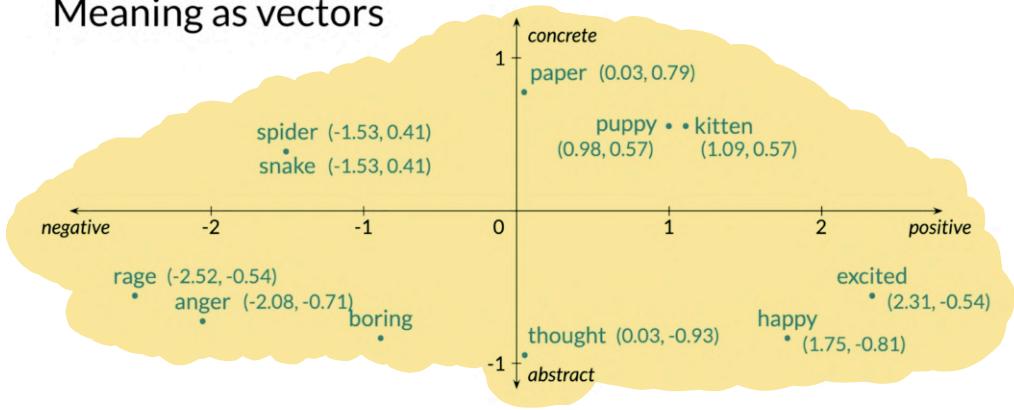
a
...
happy
...
Zyzyva



Meaning as vectors

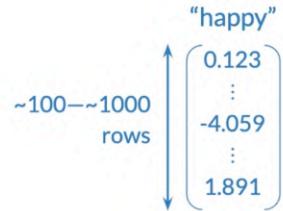


Meaning as vectors



Word embedding vectors

- + Low dimension
- + Embed meaning
 - e.g. semantic distance
forest ≈ tree forest ≠ ticket
 - e.g. analogies
Paris:France :: Rome:?



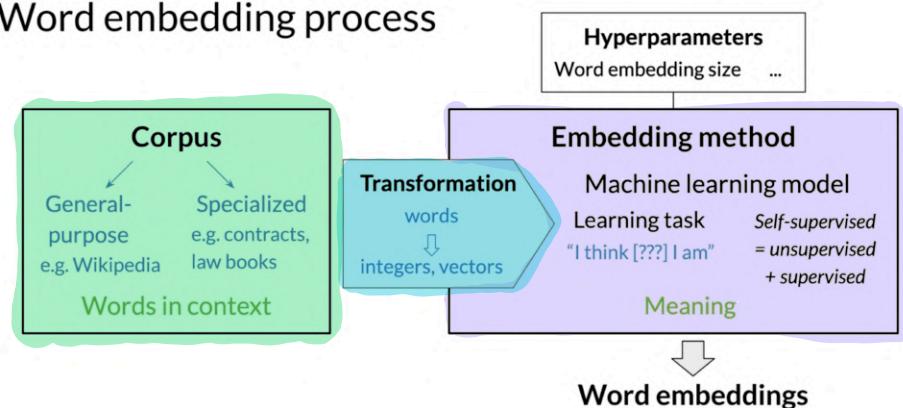
Terminology

word vectors		
integers	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

The model learns to predict the word surrounding a given input word.

Global vectors or glove by Stanford, which involves factorizing the logarithm of the corpuses word co-occurrence matrix, which is similar to the counter matrix you've used before.

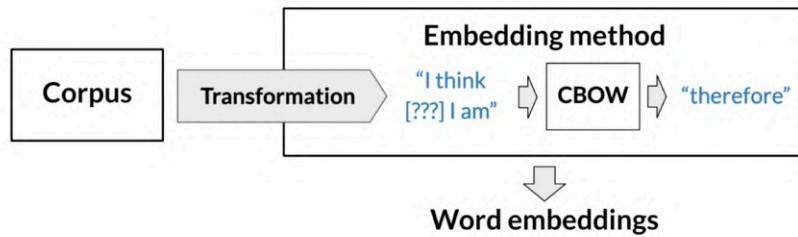
fastText by Facebook, which is based on the skip-gram model and takes into account the structure of words by representing words as an n-gram of characters. This enables the model to support previously unseen words, known as outer vocabulary words, by inferring their embedding from the sequence of characters they are made of and the corresponding sequences that it was initially trained on. For example, it would create similar embeddings for kitty and kitten, even if it had never seen the word kitty before. As kitty and kitten are made of similar sequences of characters. Another benefit of fastText, is that word embedding vectors can be averaged together to make vector representations of phrases and sentences.

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

As a reminder, to create word embeddings, you need the corpus and the machine learning model that performs a learning task. One byproduct of the learning task is a set of word embeddings. You also need a way to transform the corpus into a representation that is suited to the machine learning model. In the case of the continuous bag-of-words model, the objective of the task is to predict a missing word based on the surrounding words. The rationale is that if two unique words are both frequently surrounded by a similar sets of words when used in various sentences, then those two words tend to be related in their meaning. Another way to say this is that they are related semantically. For example, in the sentence, the little something is barking. With a large enough corpus, the model will learn to predict that the missing word is related to dogs, such as the word dog itself or puppy, hound, terrier, and so on. So the model will end up learning the meaning of words based on their context.

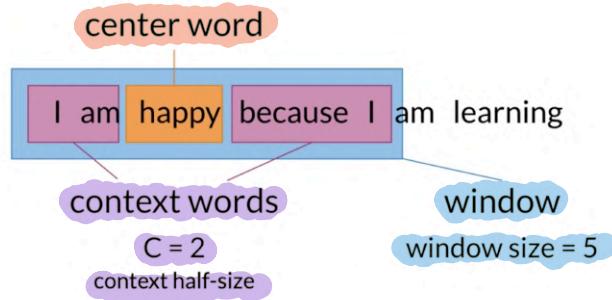
The little _____ ? is barking



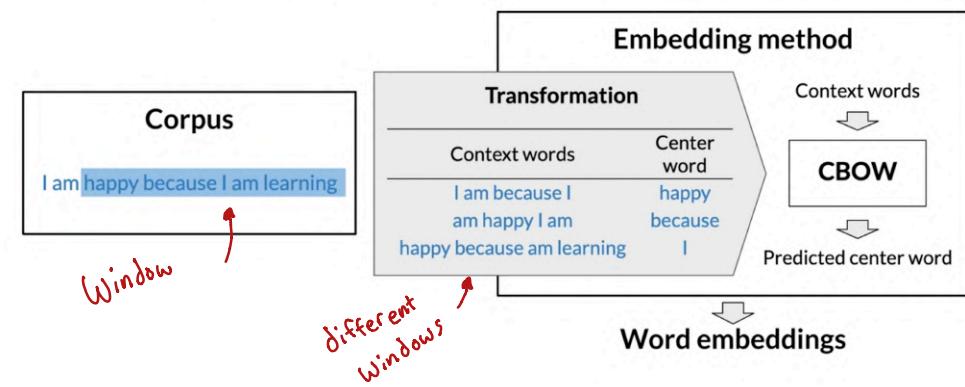
dog
puppy
hound
terrier
...



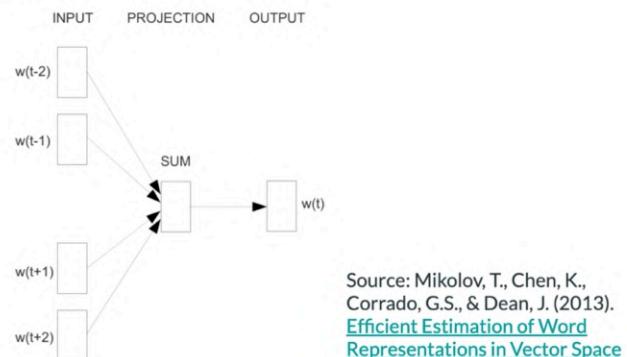
Creating a training example



From corpus to training



CBOW in a nutshell



Cleaning and tokenization matters

- Letter case "The" == "the" == "THE" → lowercase / uppercase
- 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

number of context words

The context have size stored in the variable C, which is the number of words to be taken on each side of the center word

```
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
```

I	am	happy	because	I	am	learning
0	1	2	3	4	5	6

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
```

I	am	happy	because	I	am	learning
0	1	2	3	4	5	6

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
```

I	am	happy	because	I	am	learning
0	1	2	3	4	5	6

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
```

I	am	happy	because	I	am	learning
0	1	2	3	4	5	6

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 { Vocabulary { One-hot vector	I am happy because I am learning am, because, happy, I, learning <table border="1" style="border-collapse: collapse; text-align: center;"> <tr><td>am</td><td>1</td></tr> <tr><td>because</td><td>0</td></tr> <tr><td>happy</td><td>0</td></tr> <tr><td>I</td><td>0</td></tr> <tr><td>learning</td><td>0</td></tr> </table> <table border="1" style="border-collapse: collapse; text-align: center;"> <tr><td>0</td></tr> <tr><td>1</td></tr> <tr><td>0</td></tr> <tr><td>0</td></tr> <tr><td>0</td></tr> </table> <table border="1" style="border-collapse: collapse; text-align: center;"> <tr><td>0</td></tr> <tr><td>0</td></tr> <tr><td>1</td></tr> <tr><td>0</td></tr> <tr><td>0</td></tr> </table> <table border="1" style="border-collapse: collapse; text-align: center;"> <tr><td>0</td></tr> <tr><td>0</td></tr> <tr><td>0</td></tr> <tr><td>1</td></tr> <tr><td>0</td></tr> </table> <table border="1" style="border-collapse: collapse; text-align: center;"> <tr><td>0</td></tr> <tr><td>0</td></tr> <tr><td>0</td></tr> <tr><td>0</td></tr> <tr><td>1</td></tr> </table>	am	1	because	0	happy	0	I	0	learning	0	0	1	0	0	0	0	0	1	0	0	0	0	0	1	0	0	0	0	0	1
am	1																														
because	0																														
happy	0																														
I	0																														
learning	0																														
0																															
1																															
0																															
0																															
0																															
0																															
0																															
1																															
0																															
0																															
0																															
0																															
0																															
1																															
0																															
0																															
0																															
0																															
0																															
1																															

Transforming context words into vectors

Average of individual one-hot vectors

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

Final prepared training set

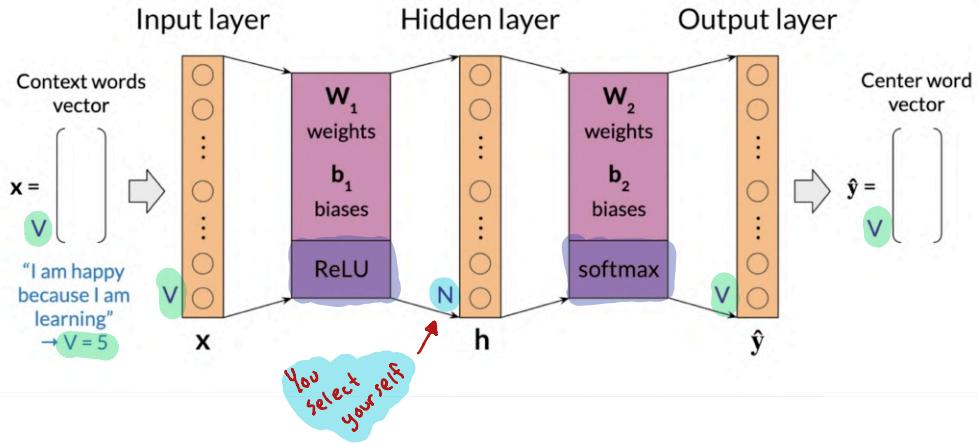
Context words	Context words vector	Center word	Center word vector
<u>I am because I</u>	[0.25; 0.25; 0; 0.5; 0]	<u>happy</u>	[0; 0; 1; 0; 0]
<u>am happy I am</u>	[0.5; 0; 0.25; 0.25; 0]	<u>because</u>	[0; 1; 0; 0; 0]
<u>happy because am learning</u>	[0.25; 0.25; 0.25; 0; 0.25]	<u>I</u>	[0; 0; 0; 1; 0]

Sliding window ↗

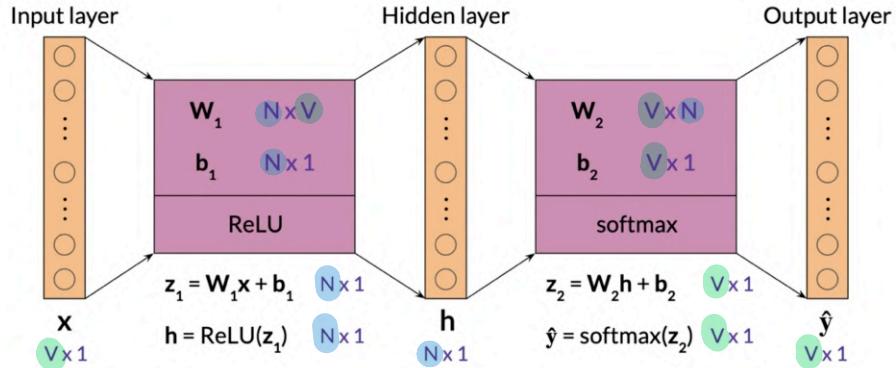
Notice, Context word has vector where center word would be

Architecture of the CBOW model

Hyperparameters
N: Word embedding size ...



Dimensions (single input)



Dimensions (single input)

Column vectors

$$z_1 = W_1 x + b_1 \quad z_1 = \begin{bmatrix} \vdots \\ N \times 1 \end{bmatrix} \quad W_1 = \begin{bmatrix} \quad N \times V \quad \end{bmatrix} \quad x = \begin{bmatrix} \vdots \\ V \times 1 \end{bmatrix} \quad b_1 = \begin{bmatrix} \vdots \\ N \times 1 \end{bmatrix}$$

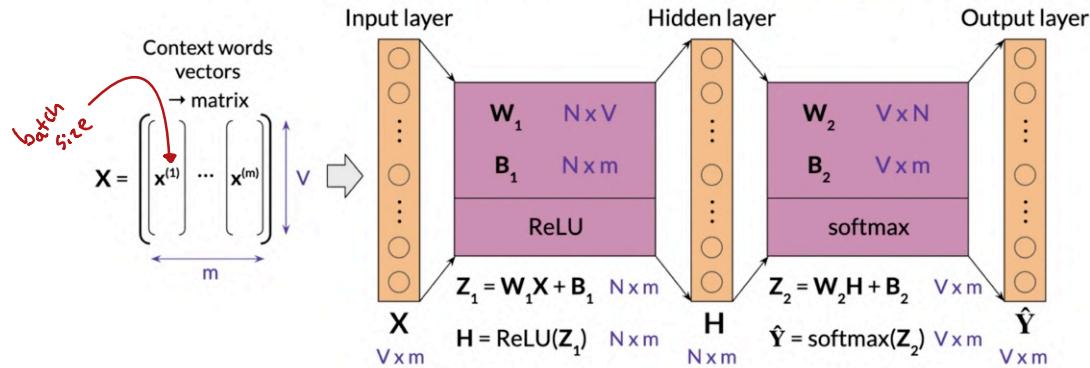
Row vectors

$$z_1 = x W_1^T + b_1 \quad b_1 = \begin{bmatrix} 1 \times N \end{bmatrix} \quad W_1 = \begin{bmatrix} \quad N \times V \quad \end{bmatrix} \quad b_1 = \begin{bmatrix} 1 \times N \end{bmatrix}$$

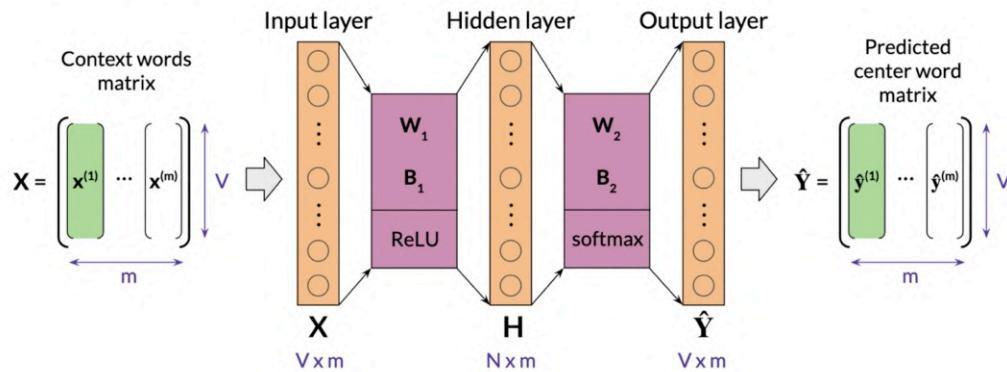
$$x = \begin{bmatrix} 1 \times V \end{bmatrix}$$

Dimensions (batch input)

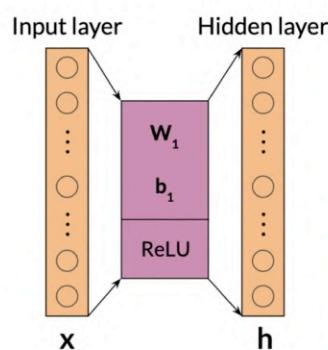
$$\begin{bmatrix} \mathbf{b}_1 \\ \vdots \\ \mathbf{b}_m \end{bmatrix} \rightarrow \mathbf{B}_1 = \begin{bmatrix} \mathbf{b}_1 & \dots & \mathbf{b}_1 \end{bmatrix} \xrightarrow{N \text{ broadcasting}}$$



Dimensions (batch input)



Rectified Linear Unit (ReLU)



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

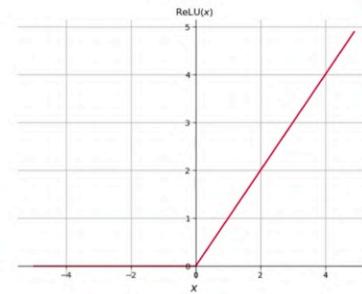
$$h = \text{ReLU}(z_1)$$

$$z_1 = \begin{bmatrix} 5.1 \\ -0.3 \\ \vdots \\ -4.6 \\ 0.2 \end{bmatrix}$$

$$\xrightarrow{\text{ReLU}}$$

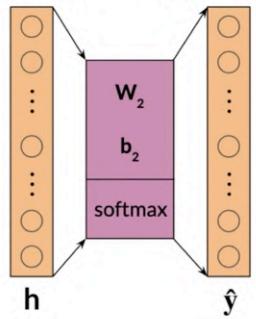
$$h = \begin{bmatrix} 5.1 \\ 0 \\ \vdots \\ 0 \\ 0.2 \end{bmatrix}$$

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



Softmax

Hidden layer Output layer



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

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

$$\left[\begin{array}{c} \end{array} \right] \in \mathbb{R}^K \xrightarrow{\text{softmax}} \left[\begin{array}{c} \end{array} \right] \in [0, 1]^K$$

\downarrow
 $\Sigma=1$

~probabilities

$$\begin{aligned} z &= \begin{pmatrix} z_1 \\ z_2 \\ \vdots \\ z_i \\ \vdots \\ z_V \end{pmatrix} \\ \hat{\mathbf{y}} &= \begin{pmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_i \\ \vdots \\ \hat{y}_V \end{pmatrix} \end{aligned}$$

$\xrightarrow{\text{softmax}}$

Probabilities of being center word

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

Softmax: example

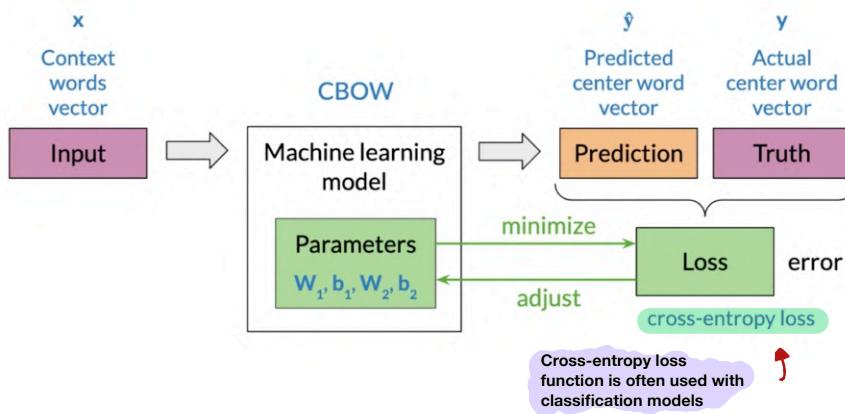
$$\hat{y}_i = \frac{e^{z_i}}{\sum_{j=1}^V e^{z_j}}$$

$$\begin{aligned} z &= \begin{pmatrix} 9 \\ 8 \\ 11 \\ 10 \\ 8.5 \end{pmatrix} \\ \exp(z) &= \begin{pmatrix} 8103 \\ 2981 \\ 59874 \\ 22026 \\ 4915 \end{pmatrix} \\ \hat{\mathbf{y}} = \text{softmax}(z) &= \begin{pmatrix} 0.083 \\ 0.03 \\ 0.612 \\ 0.225 \\ 0.05 \end{pmatrix} \end{aligned}$$

← Predicted center word

$\Sigma=1$

Loss

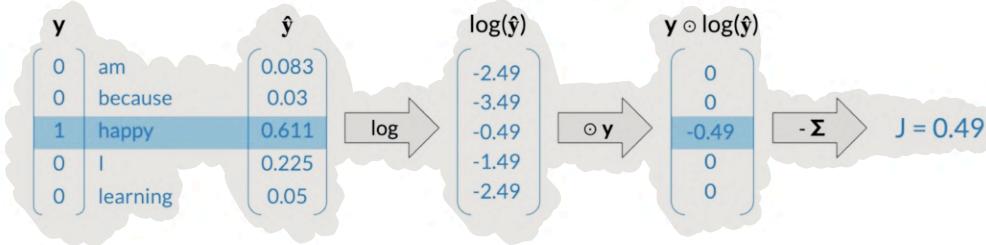


Cross-entropy loss

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

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

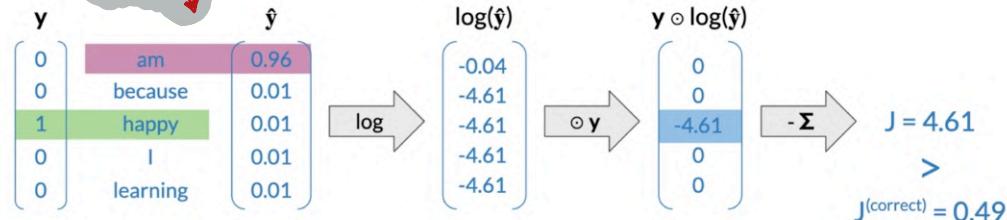
I am happy because I am learning



Cross-entropy loss

What if prediction is wrong?

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



Cross-entropy loss function is often used with classification models

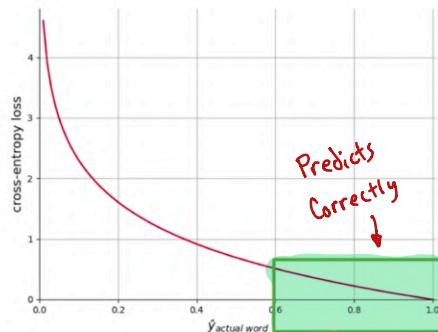
Cross-entropy loss

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

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

\mathbf{y}	$\hat{\mathbf{y}}$
$\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$	$\begin{bmatrix} \text{am} & 0.96 \\ \text{because} & 0.01 \\ \text{happy} & 0.01 \\ \text{I} & 0.01 \\ \text{learning} & 0.01 \end{bmatrix}$

$\rightarrow J = 4.61$



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	

$\rightarrow J = 4.61$

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



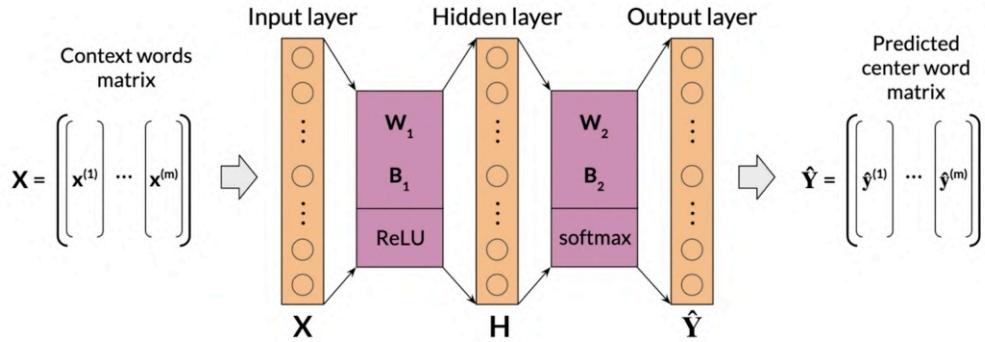
Training process

- Forward propagation
- Cost
- Backpropagation and gradient descent

Forward propagation

$$Z_1 = \mathbf{W}_1 \mathbf{X} + \mathbf{B}_1 \quad Z_2 = \mathbf{W}_2 \mathbf{H} + \mathbf{B}_2$$

$$\mathbf{H} = \text{ReLU}(Z_1) \quad \hat{\mathbf{Y}} = \text{softmax}(Z_2)$$



Cost

Cost: mean of losses

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

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

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

Predicted center word matrix: $\hat{\mathbf{Y}} = \begin{pmatrix} \hat{\mathbf{y}}^{(1)} \\ \vdots \\ \hat{\mathbf{y}}^{(m)} \end{pmatrix}$

Actual center word matrix: $\mathbf{Y} = \begin{pmatrix} \mathbf{y}^{(1)} \\ \vdots \\ \mathbf{y}^{(m)} \end{pmatrix}$

Minimizing the cost

$$J_{batch} = f(\mathbf{W}_1, \mathbf{W}_2, \mathbf{b}_1, \mathbf{b}_2)$$

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

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

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

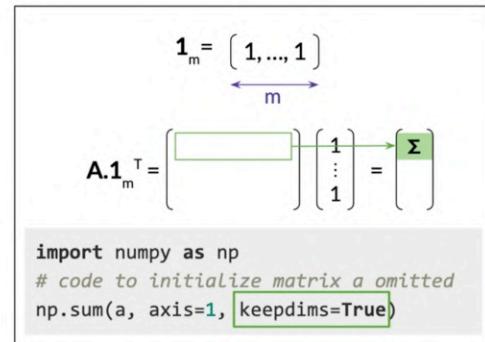
Backpropagation

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

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

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

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



Sum of
columns

Set the `keepdims` parameter to True, so that the results can be broadcasts into a matrix of whatever size is necessary later in the calculations.

Gradient descent

Hyperparameter: learning rate α

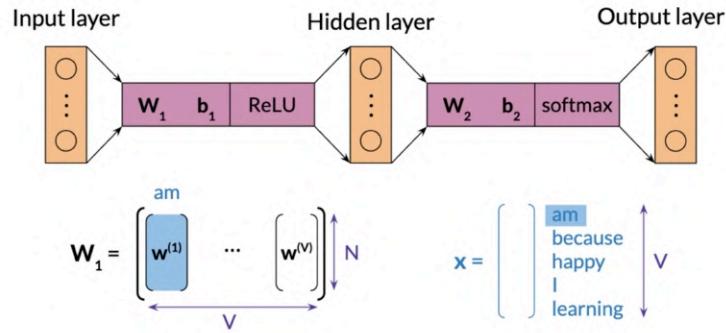
$$\mathbf{W}_1 := \mathbf{W}_1 - \alpha \frac{\partial J_{batch}}{\partial \mathbf{W}_1}$$

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

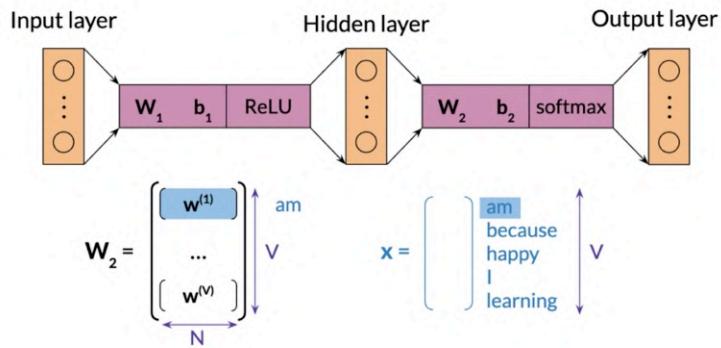
$$\mathbf{b}_1 := \mathbf{b}_1 - \alpha \frac{\partial J_{batch}}{\partial \mathbf{b}_1}$$

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

Extracting word embedding vectors: option 1



Extracting word embedding vectors: option 2



Extracting word embedding vectors: option 3

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

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

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

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

Intrinsic evaluation

Test relationships between words

- ## • Analogies

Semantic analogies

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

You should be aware of
for this approach, is that
there could be several
possible correct
answers

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

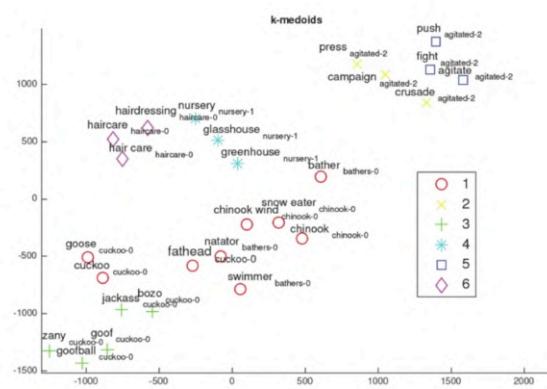
Intrinsic evaluation

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)
```

