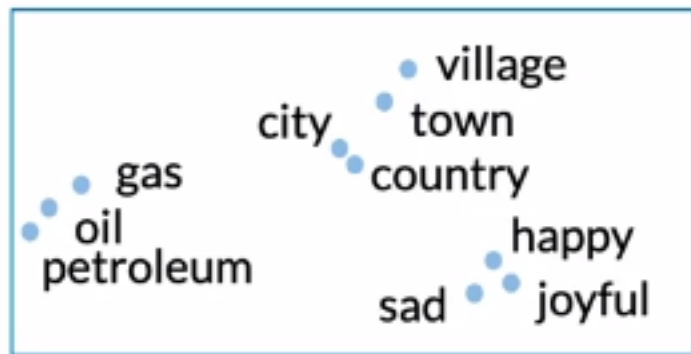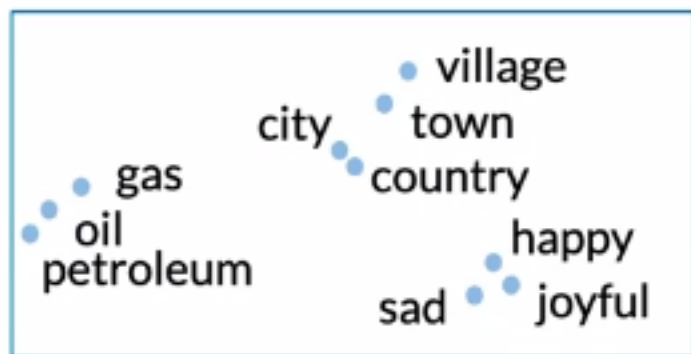# Some basic applications of word embeddings

# Some basic applications of word embeddings



Semantic analogies
and similarity

# Some basic applications of word embeddings



Semantic analogies
and similarity



Sentiment analysis

# Some basic applications of word embeddings



Semantic analogies
and similarity



Sentiment analysis



Classification of
customer feedback

# Advanced applications of word embeddings



Machine translation

# Advanced applications of word embeddings



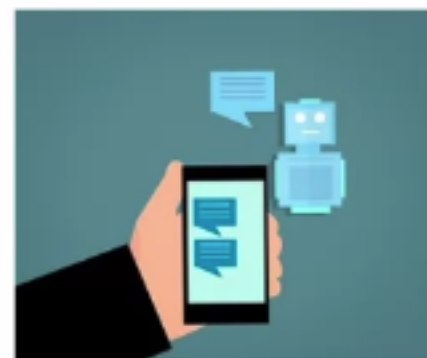Machine translation



Information extraction

# Advanced applications of word embeddings



Machine translation



Information extraction



Question answering

# Learning objectives

- Identify the key concepts of word representations

# Learning objectives

- Identify the key concepts of word representations

- Generate word embeddings

# Learning objectives

- Identify the key concepts of word representations

- Generate word embeddings

- Prepare text for machine learning

# Learning objectives

- Identify the key concepts of word representations

- Generate word embeddings

- Prepare text for machine learning

- Implement the continuous bag-of-words model

# Learning objectives

- 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

# Integers

+   Simple

-   Ordering: little semantic sense

**hand** < **happy** < **zebra**

615   ?!   621   ?!   1000

# One-hot vectors

$$
\left.
\begin{array}{l}
\text{a} \\
\text{able} \\
\text{about} \\
\ldots \\
\text{hand} \\
\ldots \\
\text{happy} \\
\ldots \\
\text{zebra}
\end{array}
\right\} \; 1000 \text{ rows}
$$

# One-hot vectors

# One-hot vectors



"a"

| 1 | a |
| 0 | able |
| 0 | about |
| ⋮ | ... |
| 0 | hand |
| ⋮ | ... |
| 0 | happy |
| ⋮ | ... |
| 0 | zebra |

•••

1000 rows

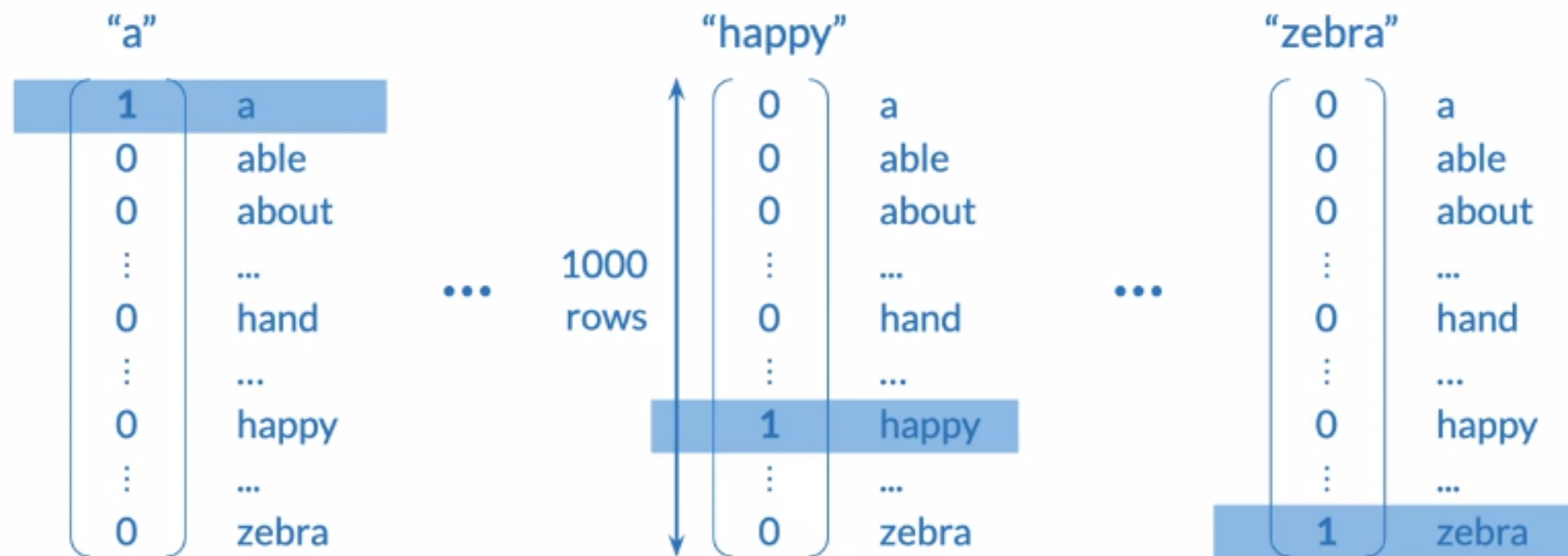"happy"

| 0 | a |
| 0 | able |
| 0 | about |
| ⋮ | ... |
| 0 | hand |
| ⋮ | ... |
| 1 | happy |
| ⋮ | ... |
| 0 | zebra |

•••

"zebra"

| 0 | a |
| 0 | able |
| 0 | about |
| ⋮ | ... |
| 0 | hand |
| ⋮ | ... |
| 0 | happy |
| ⋮ | ... |
| 1 | zebra |

# One-hot vectors

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

"happy"

| | | |
|---|---|---|
| 1 | 0 | a |
| 2 | 0 | able |
| 3 | 0 | about |
| ... | ⋮ | ... |
| 615 | 0 | hand |
| ... | ⋮ | ... |
| 621 | 1 | happy |
| ... | ⋮ | ... |
| 1000 | 0 | zebra |

# One-hot vectors

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

"happy"

$$
\begin{array}{ccl}
1 & 0 & a \\
2 & 0 & able \\
3 & 0 & about \\
... & \vdots & ... \\
615 & 0 & hand \\
... & \vdots & ... \\
621 & 1 & happy \\
... & \vdots & ... \\
1000 & 0 & zebra
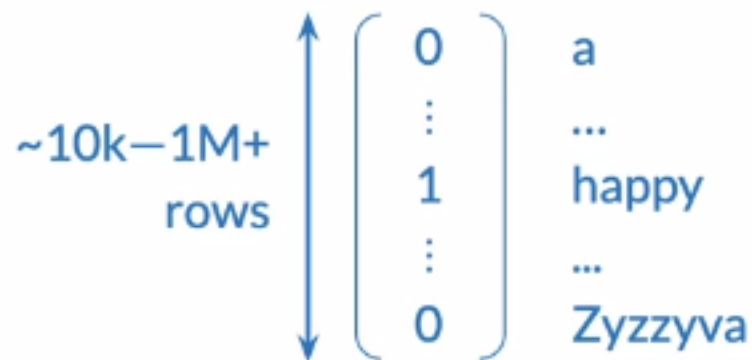\end{array}
$$

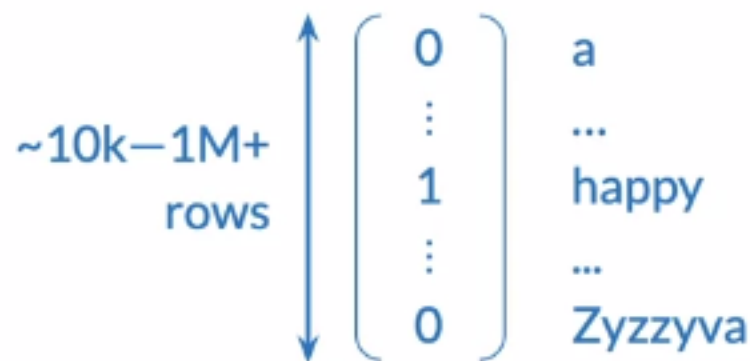621 → 621

# One-hot vectors

   +   Simple

   +   No implied ordering

# One-hot vectors

+ Simple

+ No implied ordering

- Huge vectors

# One-hot vectors
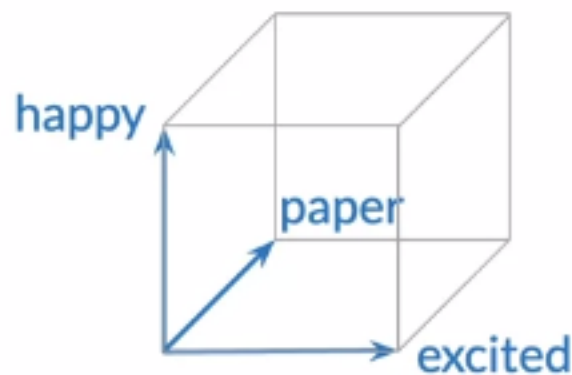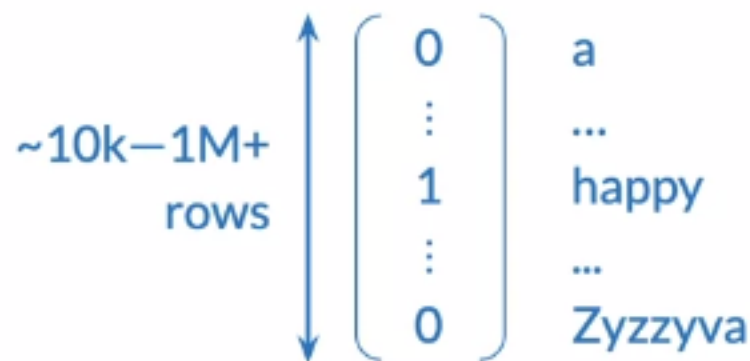
+  Simple

+  No implied ordering

-  Huge vectors

$$\sim 10k{-}1M+ \text{ rows} \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \begin{matrix} a \\ \ldots \\ happy \\ \ldots \\ Zyzzyva \end{matrix}$$

# One-hot vectors

+ Simple

+ No implied ordering

- Huge vectors

- No embedded meaning

$$\sim 10k-1M+ \text{ rows} \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \begin{matrix} a \\ \dots \\ happy \\ \dots \\ Zyzzyva \end{matrix}$$

# One-hot vectors

+ Simple

+ No implied ordering

- Huge vectors

- No embedded meaning

$$\sim 10k - 1M+ \text{ rows} \begin{pmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{pmatrix} \begin{matrix} a \\ \dots \\ happy \\ \dots \\ Zyzzyva \end{matrix}$$

happy

paper

excited

d(paper, excited)
= d(paper, happy)
= d(excited, happy)

deeplearning.ai

Word
Embeddings

# Meaning as vectors



negative     -2     -1     0     1     2     positive

# Meaning as vectors

# Meaning as vectors



rage (-2.52)  anger (-2.08)  spider (-1.53)  boring (-0.91)  paper (0.03)  kitten (1.09)  happy (1.75)  excited (2.31)

negative  -2  -1  0  1  2  positive
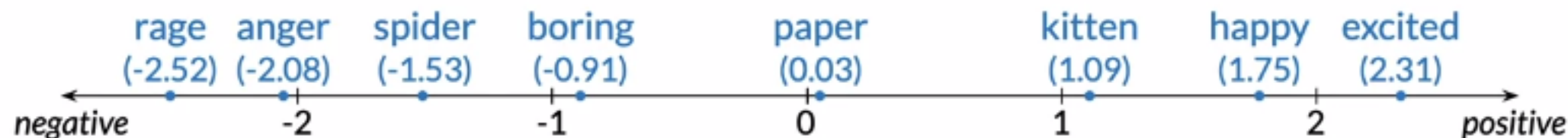
# Meaning as vectors

# Meaning as vectors

Meaning as vectors

concrete

1 — paper

spider
snake

negative   -2        -1        0                1                2   positive

rage
anger
boring          thought
                -1
happy   excited
abstract

puppy ● ● kitten

deeplearning.ai

# Meaning as vectors



concrete

1 ┤

paper  (0.03, 0.79)

spider  (-1.53, 0.41)

snake  (-1.53, 0.41)

puppy • • kitten
(0.98, 0.57)      (1.09, 0.57)

negative ← -2 ─── -1 ─── 0 ─── 1 ─── 2 → positive

rage  (-2.52, -0.54)

anger  (-2.08, -0.71)

boring

thought  (0.03, -0.93)

excited
(2.31, -0.54)

happy
(1.75, -0.81)

-1 ┤
abstract

deeplearning.ai

# Word embedding vectors

# Word embedding vectors

+ Low dimension

# Word embedding vectors

+ Low dimension

$$
\text{~100—~1000 rows} \left\{ \begin{array}{c} \text{"happy"} \\ 0.123 \\ \vdots \\ -4.059 \\ \vdots \\ 1.891 \end{array} \right.
$$

# Word embedding vectors

+ Low dimension

+ Embed meaning

"happy"

~100—~1000 rows

$$\begin{pmatrix} 0.123 \\ \vdots \\ -4.059 \\ \vdots \\ 1.891 \end{pmatrix}$$

# Word embedding vectors

+ **Low dimension**

+ **Embed meaning**
  - **e.g. semantic distance**

"happy"

$\sim$100—$\sim$1000 rows

$$\begin{pmatrix} 0.123 \\ \vdots \\ -4.059 \\ \vdots \\ 1.891 \end{pmatrix}$$

# Word embedding vectors

+ **Low dimension**

+ **Embed meaning**
  ○ **e.g. semantic distance**

forest ≈ tree     forest ≠ ticket

"happy"

~100—~1000
rows

$\begin{pmatrix} 0.123 \\ \vdots \\ -4.059 \\ \vdots \\ 1.891 \end{pmatrix}$

# Word embedding vectors

+ **Low dimension**

+ **Embed meaning**
  - ○ **e.g. semantic distance**

    forest ≈ tree      forest ≠ ticket

  - ○ **e.g. analogies**

"happy"

$$\sim100\text{—}\sim1000 \text{ rows} \begin{pmatrix} 0.123 \\ \vdots \\ -4.059 \\ \vdots \\ 1.891 \end{pmatrix}$$
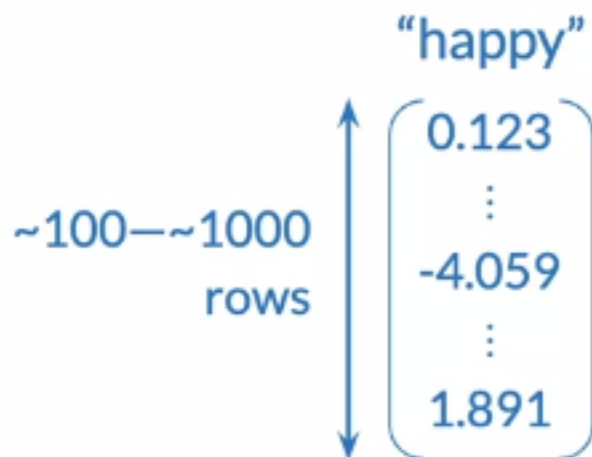
# Word embedding vectors

+ Low dimension

+ Embed meaning
  ○ e.g. semantic distance

  forest ≈ tree     forest ≠ ticket

  ○ e.g. analogies

  Paris:France :: Rome:?

"happy"

$\sim$100—$\sim$1000 rows

$$\begin{pmatrix} 0.123 \\ \vdots \\ -4.059 \\ \vdots \\ 1.891 \end{pmatrix}$$

# Terminology

integers          one-hot vectors          word embedding vectors

# Terminology

integers

word vectors

one-hot vectors    word embedding vectors

# 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

| Corpus | Embedding method |
|--------|------------------|
|        |                  |

# Word embedding process

| Corpus | Embedding method |
|---|---|
| **Words in context** | |

# Word embedding process

Corpus

General-
purpose
e.g. Wikipedia

Words in context

# Word embedding process



Corpus

General-purpose → Specialized

e.g. Wikipedia

Words in context

# Word embedding process



**Corpus**

General-purpose
e.g. Wikipedia

Specialized
e.g. contracts,
law books

Words in context

# Word embedding process

**Corpus**

General-purpose
e.g. Wikipedia

Specialized
e.g. contracts,
law books

Words in context

**Embedding method**

# Word embedding process

**Corpus**

General-purpose
e.g. Wikipedia

Specialized
e.g. contracts, law books

Words in context

**Embedding method**

Machine learning model

# Word embedding process



**Corpus**

General-purpose
e.g. Wikipedia

Specialized
e.g. contracts,
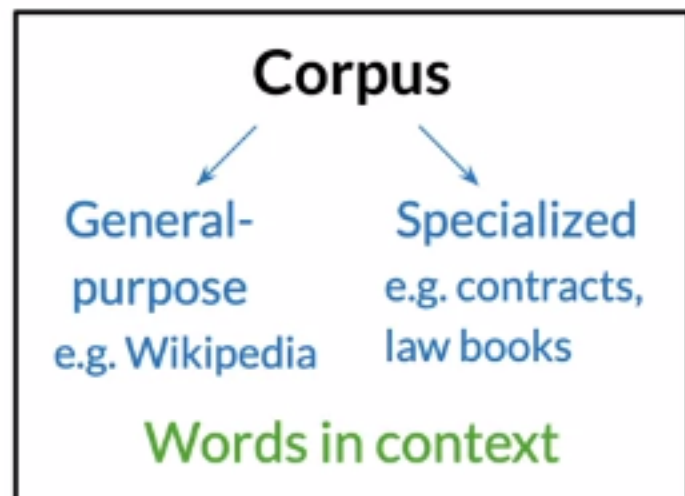law books

Words in context

**Embedding method**

Machine learning model

Learning task

Word embeddings

# Word embedding process



**Corpus**

General-purpose
e.g. Wikipedia

Specialized
e.g. contracts,
law books

Words in context

**Embedding method**

Machine learning model

Learning task

"I think [???] I am"

⬇

**Word embeddings**

# Word embedding process

**Corpus**

General-purpose
e.g. Wikipedia

Specialized
e.g. contracts,
law books

Words in context

**Embedding method**

Machine learning model

Learning task

"I think [???] I am"

Meaning

**Word embeddings**

# Word embedding process

**Corpus**

General-purpose
e.g. Wikipedia

Specialized
e.g. contracts, law books

Words in context

**Embedding method**

Machine learning model

Learning task          *Self-supervised*
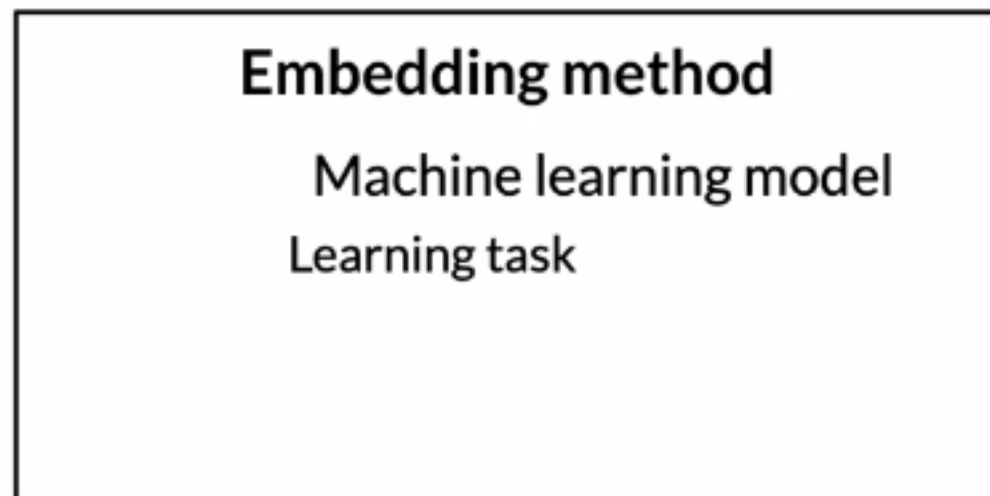
"I think [???] I am"    *= unsupervised*
                        *+ supervised*

Meaning

**Word embeddings**

# Word embedding process

**Hyperparameters**
Word embedding size     ...

**Corpus**

General-purpose
e.g. Wikipedia

Specialized
e.g. contracts,
law books

Words in context

**Transformation**

**Embedding method**

Machine learning model
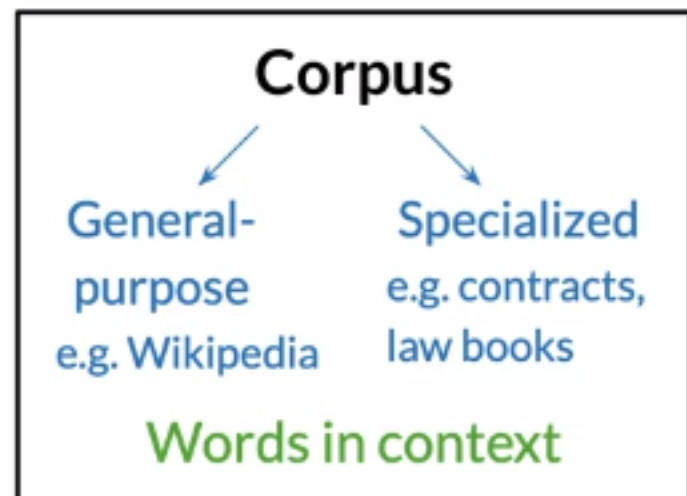
Learning task
"I think [???] I am"

*Self-supervised
= unsupervised
+ supervised*

Meaning

**Word embeddings**

# Word embedding process

**Hyperparameters**
Word embedding size     ...

**Corpus**

General-purpose
e.g. Wikipedia

Specialized
e.g. contracts, law books

Words in context

**Transformation**

words

⇩

integers, vectors

**Embedding method**

Machine learning model

Learning task

"I think [???] I am"

*Self-supervised*
*= unsupervised*
*+ supervised*

Meaning

**Word embeddings**

# Word embedding process

**Hyperparameters**
Word embedding size    ...

**Corpus**

General-purpose
e.g. Wikipedia

Specialized
e.g. contracts,
law books

Words in context

**Transformation**
words
⇩
integers, vectors

**Embedding method**

Machine learning model
Learning task
"I think [???] I am"

*Self-supervised
= unsupervised
+ supervised*

Meaning

**Word embeddings**

deeplearning.ai

# Basic word embedding methods

- word2vec (Google, 2013)
    - Continuous bag-of-words (CBOW)

# Basic word embedding methods

- word2vec (Google, 2013)
    - Continuous bag-of-words (CBOW)
    - Continuous skip-gram / Skip-gram with negative sampling (SGNS)

# 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)

# 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

# Advanced word embedding methods

Deep learning, contextual embeddings

- BERT (Google, 2018)

# Advanced word embedding methods

Deep learning, contextual embeddings

- BERT (Google, 2018)

- ELMo (Allen Institute for AI, 2018)

# Advanced word embedding methods

Deep learning, contextual embeddings

- BERT (Google, 2018)

- ELMo (Allen Institute for AI, 2018)

- GPT-2 (OpenAI, 2018)

# 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

# Continuous bag-of-words word embedding process

Corpus

Embedding method

# Continuous bag-of-words word embedding process

Corpus

Embedding method

Word embeddings

# Continuous bag-of-words word embedding process



Corpus → Transformation → **Embedding method** → **Word embeddings**

# Continuous bag-of-words word embedding process

# Center word prediction: rationale

# Center word prediction: rationale

The little ____?____ is barking

# Creating a training example

I am happy because I am learning

# Creating a training example

center word

I  am  happy  because  I  am  learning

# Creating a training example

center word

I am **happy** because I am learning

context words

# From corpus to training

Corpus → Transformation → CBOW

Embedding method

Corpus

Transformation → Context words → CBOW → Predicted center word

Word embeddings

# From corpus to training

**Embedding method**

**Corpus**

**Transformation**

Context words                Center word

Context words

**CBOW**

Predicted center word

**Word embeddings**

deeplearning.ai

# From corpus to training

**Corpus**

I am happy because I am learning

**Embedding method**

**Transformation**

| Context words | Center word |
|---|---|
| I am because I | happy |

Context words

**CBOW**

Predicted center word

**Word embeddings**

deeplearning.ai

# CBOW in a nutshell



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

deeplearning.ai

# Cleaning and tokenization matters

# Cleaning and tokenization matters

- Letter case

# Cleaning and tokenization matters

- Letter case          "The" == "the" == "THE"

# Cleaning and tokenization matters

- Letter case      "The" == "the" == "THE"    → *lowercase / upper case*

# Cleaning and tokenization matters

- Letter case        "The" == "the" == "THE"    → *lowercase / upper case*

- Punctuation

# Cleaning and tokenization matters

- Letter case      "The" == "the" == "THE"    → *lowercase / upper case*

- Punctuation      , ! . ? → .

# Cleaning and tokenization matters

- Letter case       "The" == "the" == "THE"    → *lowercase / upper case*

- Punctuation       , ! . ? → .       " ' « » ' " → ∅

# Cleaning and tokenization matters

- Letter case        "The" == "the" == "THE"    → *lowercase / upper case*

- Punctuation       , ! . ? → .      " ' « » ' " → ∅      ... !! ??? → .

# Cleaning and tokenization matters

- Letter case      "The" == "the" == "THE"   → *lowercase / upper case*

- Punctuation      , ! . ? → .        " ' « » ' " → ∅        ... !! ??? → .

- Numbers

# Cleaning and tokenization matters

- Letter case      "The" == "the" == "THE"    → *lowercase / upper case*

- Punctuation      , ! . ? → .      " ' « » ' " → ∅      ... !! ??? → .

- Numbers      1 2 3 5 8 → ∅      3.14159   90210

# Cleaning and tokenization matters

- Letter case                "The" == "the" == "THE"   → *lowercase / upper case*

- Punctuation           , ! . ? → .           " ' « » ' " → ∅        … !! ??? → .

- Numbers             1 2 3 5 8 → ∅       3.14159  90210

# 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>*

# 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

# 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       ∇ $ € § ¶ **

# 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      ∇ $ € § ¶ **     → ∅

# 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

# 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

# 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!!!

# Example in Python: corpus

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

| emoji | punctuation | number |

# Example in Python: libraries

```python
# 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

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

data = re.sub(r'[,!?;-]+', '.', corpus)
```

# Example in Python: code

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

data = re.sub(r'[,!?;-]+', '.', corpus)
```

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

# Example in Python: code

```python
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

```python
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)
       ]
```

# Example in Python: code

```python
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

```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

```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

```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

```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 |
|---|----|----|----|----|----|----|

# Sliding window of words in Python

```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

```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

```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

```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

```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

```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
```

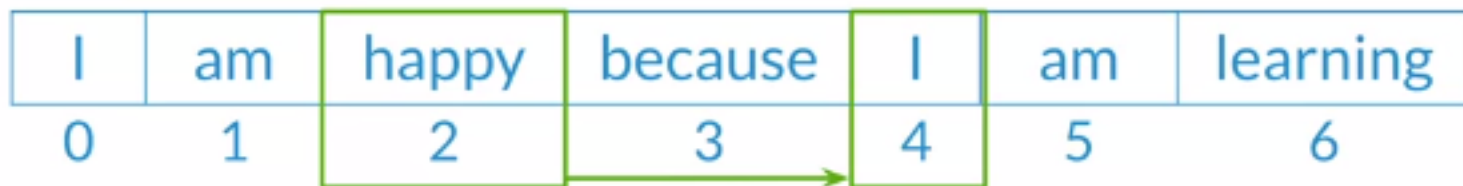| I | am | happy | because | I | am | learning |
|---|----|-------|---------|---|----|----------|
| 0 | 1  | 2     | 3       | 4 | 5  | 6        |

# Sliding window of words in Python

```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

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

```python
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

```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

```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

# Transforming center words into vectors

Corpus          I am happy because I am learning

Vocabulary      am, because, happy, I, learning

# Transforming center words into vectors

Corpus      I am happy because I am learning

Vocabulary      am, because, happy, I, learning

One-hot vector

$$
\begin{array}{c}
\text{am} \\
\text{because} \\
\text{happy} \\
\text{I} \\
\text{learning}
\end{array}
\begin{pmatrix}
\phantom{x} \\
\phantom{x} \\
\phantom{x} \\
\phantom{x} \\
\phantom{x}
\end{pmatrix}
$$

# Transforming center words into vectors

Corpus        I am happy because I am learning

Vocabulary    am, because, happy, I, learning

One-hot
vector

|          | am |
|----------|----|
| am       | 1  |
| because  | 0  |
| happy    | 0  |
| I        | 0  |
| learning | 0  |

# Transforming center words into vectors

Corpus      I am happy because I am learning

Vocabulary      am, because, happy, I, learning

One-hot vector

|          | am | because |
|----------|----|---------|
| am       | 1  | 0       |
| because  | 0  | 1       |
| happy    | 0  | 0       |
| I        | 0  | 0       |
| learning | 0  | 0       |

# 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       | 1  | 0       | 0     | 0 | 0        |
| 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

# Transforming context words into vectors

Average of individual one-hot vectors

I      am      because      I

# Transforming context words into vectors

Average of individual one-hot vectors

$$
\begin{array}{cccc}
\text{I} & \text{am} & \text{because} & \text{I}
\end{array}
$$

$$
\begin{array}{c}
\text{am} \\
\text{because} \\
\text{happy} \\
\text{I} \\
\text{learning}
\end{array}
\left[
\begin{array}{c}
0 \\
0 \\
0 \\
1 \\
0
\end{array}
\right]
$$

# Transforming context words into vectors

Average of individual one-hot vectors

|            | I | am | because | I |
|-----------:|:-:|:--:|:-------:|:-:|
| am         | 0 |    |         |   |
| because    | 0 |    |         |   |
| happy      | 0 |    |         |   |
| I          | 1 |    |         |   |
| learning   | 0 |    |         |   |

# Transforming context words into vectors

Average of individual one-hot vectors

$$
\begin{array}{c}
\quad\quad\text{I}\quad\quad\quad\quad\quad\text{am}\quad\quad\quad\text{because}\quad\quad\quad\text{I} \\
\begin{array}{r}
\text{am} \\
\text{because} \\
\text{happy} \\
\text{I} \\
\text{learning}
\end{array}
\left[\begin{array}{c}
0 \\
0 \\
0 \\
1 \\
0
\end{array}\right]
\left[\begin{array}{c}
1 \\
0 \\
0 \\
0 \\
0
\end{array}\right]
\end{array}
$$

# Transforming context words into vectors

Average of individual one-hot vectors

|          | I | am | because | I |
|----------|---|----|---------|---|
| am       | 0 | 1  | 0       | 0 |
| because  | 0 | 0  | 1       | 0 |
| happy    | 0 | 0  | 0       | 0 |
| I        | 1 | 0  | 0       | 1 |
| learning | 0 | 0  | 0       | 0 |

# Transforming context words into vectors

Average of individual one-hot vectors

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

deeplearning.ai

# Transforming context words into vectors

Average of individual one-hot vectors

$$
\begin{array}{c} \text{am} \\ \text{because} \\ \text{happy} \\ \text{I} \\ \text{learning} \end{array}
\begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}^{\text{I}}
+
\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}^{\text{am}}
+
\begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}^{\text{because}}
+
\begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}^{\text{I}}
\;/\,4\;=\;
\begin{bmatrix} 0.25 \\ 0.25 \\ 0 \\ 0.5 \\ 0 \end{bmatrix}^{\text{I am because I}}
$$

# 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] |

# 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

# Architecture of the CBOW model

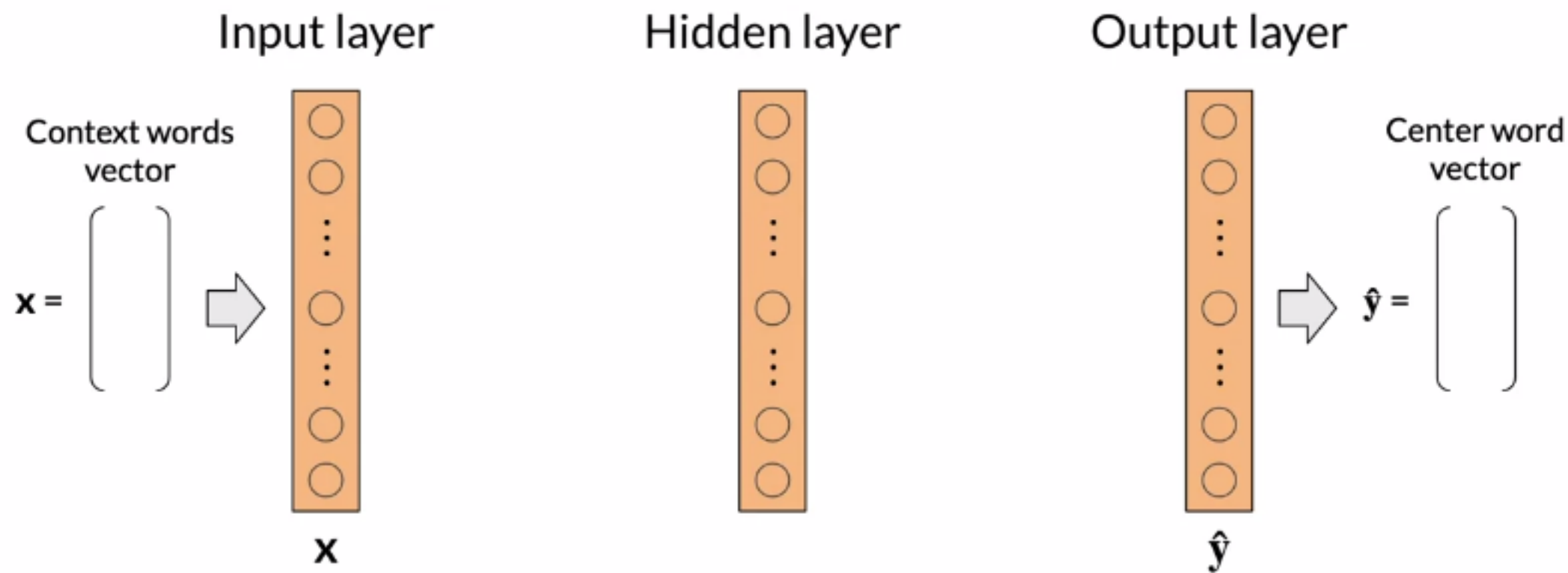Input layer       Hidden layer       Output layer

# Architecture of the CBOW model



Input layer       Hidden layer       Output layer

Context words vector

$\mathbf{x} =$

$\mathbf{x}$

Center word vector

$\hat{\mathbf{y}} =$

$\hat{\mathbf{y}}$

# Architecture of the CBOW model

Input layer        Hidden layer        Output layer



Context words vector

$$\mathbf{x} = \begin{bmatrix} \\ \\ \\ \\ \end{bmatrix} V$$

$\mathbf{x}$

Center word vector

$$\hat{\mathbf{y}} = \begin{bmatrix} \\ \\ \\ \\ \end{bmatrix} V$$

$\hat{\mathbf{y}}$

# Architecture of the CBOW model

Input layer             Hidden layer             Output layer

Context words
vector

$$\mathbf{x} =$$

V

"I am happy
because I am
learning"
→ V = 5



**x**

**ŷ**

Center word
vector

$$\hat{\mathbf{y}} =$$

V

deeplearning.ai

# Architecture of the CBOW model

Input layer

Hidden layer

Output layer

Context words
vector

$\mathbf{x} =$

V

"I am happy
because I am
learning"
$\rightarrow V = 5$

Center word
vector

$\hat{\mathbf{y}} =$

V



V

V

V

V

$\mathbf{x}$

$\hat{\mathbf{y}}$

# Architecture of the CBOW model

Input layer

Hidden layer

Output layer

Context words vector

$$\mathbf{x} = \begin{bmatrix} \\ \\ \\ \\ \end{bmatrix}$$
V

"I am happy because I am learning"
→ V = 5

V

x

N

Center word vector

$$\hat{\mathbf{y}} = \begin{bmatrix} \\ \\ \\ \\ \end{bmatrix}$$
V

V

$\hat{y}$

deeplearning.ai

# Architecture of the CBOW model

Input layer

Hidden layer

Output layer



Context words vector

$\mathbf{x} =$

V

"I am happy because I am learning"
→ V = 5

V

**x**

N

**h**

V

$\hat{\mathbf{y}} =$

Center word vector

V

V

**ŷ**

# Architecture of the CBOW model

Input layer      Hidden layer      Output layer

Context words
vector

$x =$

V

"I am happy
because I am
learning"
$\rightarrow$ V = 5

V

$x$

N

$h$

V

$\hat{y}$

Center word
vector

$\hat{y} =$

V

# Architecture of the CBOW model

Input layer          Hidden layer          Output layer

Context words
vector

$\mathbf{x} =$

V

"I am happy
because I am
learning"
→ V = 5

$\mathbf{x}$

$\mathbf{W_1}$
weights

$\mathbf{b_1}$
biases

$\mathbf{h}$

$\mathbf{W_2}$
weights

$\mathbf{b_2}$
biases

Center word
vector

$\hat{\mathbf{y}} =$

V

$\hat{\mathbf{y}}$

V          N          V

Architecture of the CBOW model

**Hyperparameters**
$N$: Word embedding size ...

Input layer    Hidden layer    Output layer

Context words vector

$\mathbf{x} =$

$V$

"I am happy because I am learning"
$\rightarrow V = 5$

$V$

$\mathbf{x}$

$\mathbf{W}_1$ weights

$\mathbf{b}_1$ biases

ReLU

$N$

$\mathbf{h}$

$\mathbf{W}_2$ weights

$\mathbf{b}_2$ biases

$V$

$\hat{\mathbf{y}}$

Center word vector

$\hat{\mathbf{y}} =$

$V$

# Architecture of the CBOW model

**Hyperparameters**
N: Word embedding size    ...

Input layer         Hidden layer         Output layer

Context words
vector

$\mathbf{x} =$

V

"I am happy
because I am
learning"
→ V = 5

$\mathbf{W_1}$
weights

$\mathbf{b_1}$
biases

ReLU

Center word
vector

$\mathbf{W_2}$
weights

$\mathbf{b_2}$
biases

softmax

$\hat{\mathbf{y}} =$

V

V     x     N     h     V     $\hat{\mathbf{y}}$

deeplearning.ai

Dimensions (single input)

# Dimensions (single input)

Input layer             Hidden layer             Output layer

$W_1$

$b_1$

ReLU

$z_1 = W_1 x + b_1$

$W_2$

$b_2$

softmax

$x$

$V \times 1$

$h$

$\hat{y}$

# Dimensions (single input)

Input layer              Hidden layer            Output layer

$W_1$

$b_1$

ReLU

$z_1 = W_1 x + b_1$

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

$x$

$V \times 1$

$W_2$

$b_2$

softmax

$h$

$\hat{y}$

deeplearning.ai

# Dimensions (single input)



Input layer

Hidden layer

Output layer

$W_1$    $N \times V$

$b_1$

ReLU

$z_1 = W_1 x + b_1$

$h = ReLU(z_1)$

$x$

$V \times 1$

$h$

$W_2$

$b_2$

softmax

$\hat{y}$

# Dimensions (single input)



**Input layer**

$\mathbf{x}$

$V \times 1$

$\mathbf{W_1}$    $N \times V$

$\mathbf{b_1}$    $N \times 1$

ReLU

$\mathbf{z_1} = \mathbf{W_1}\mathbf{x} + \mathbf{b_1}$

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

**Hidden layer**

$\mathbf{h}$

$\mathbf{W_2}$

$\mathbf{b_2}$

softmax

**Output layer**

$\hat{\mathbf{y}}$

# Dimensions (single input)



Input layer

Hidden layer

Output layer

$\mathbf{W_1}$   $N \times V$

$\mathbf{b_1}$   $N \times 1$

ReLU

$\mathbf{z_1} = \mathbf{W_1}\mathbf{x} + \mathbf{b_1}$   $N \times 1$

$\mathbf{h} = ReLU(\mathbf{z_1})$   $N \times 1$

$\mathbf{W_2}$

$\mathbf{b_2}$

softmax

$\mathbf{z_2} = \mathbf{W_2}\mathbf{h} + \mathbf{b_2}$

$\mathbf{x}$

$V \times 1$

$\mathbf{h}$

$N \times 1$

$\hat{\mathbf{y}}$

# Dimensions (single input)



**Input layer**

**Hidden layer**

**Output layer**

$$W_1 \quad N \times V$$

$$b_1 \quad N \times 1$$

ReLU

$$W_2$$

$$b_2$$

softmax

$$z_1 = W_1 x + b_1 \quad N \times 1$$

$$h = \text{ReLU}(z_1) \quad N \times 1$$

$$z_2 = W_2 h + b_2$$

$$\hat{y} = \text{softmax}(z_2)$$

**x**

$V \times 1$

**h**

$N \times 1$

$\hat{y}$

# Dimensions (single input)



Input layer

Hidden layer

Output layer

$$W_1 \quad N \times V$$

$$b_1 \quad N \times 1$$

ReLU

$$W_2 \quad V \times N$$

$$b_2$$

softmax

$$z_1 = W_1 x + b_1 \quad N \times 1$$

$$h = ReLU(z_1) \quad N \times 1$$

$$z_2 = W_2 h + b_2$$

$$\hat{y} = softmax(z_2)$$

$$\mathbf{x}$$
$$V \times 1$$

$$\mathbf{h}$$
$$N \times 1$$

$$\hat{y}$$

# Dimensions (single input)



Input layer          Hidden layer          Output layer

$W_1$   $N \times V$

$b_1$   $N \times 1$

ReLU

$W_2$   $V \times N$

$b_2$   $V \times 1$

softmax

$z_1 = W_1 x + b_1$   $N \times 1$

$h = \text{ReLU}(z_1)$   $N \times 1$

$z_2 = W_2 h + b_2$

$\hat{y} = \text{softmax}(z_2)$

$x$

$V \times 1$

$h$

$N \times 1$

$\hat{y}$

# Dimensions (single input)



Input layer

Hidden layer

Output layer

$$\mathbf{W}_1 \quad N \times V$$

$$\mathbf{b}_1 \quad N \times 1$$

ReLU

$$\mathbf{W}_2 \quad V \times N$$

$$\mathbf{b}_2 \quad V \times 1$$

softmax

$$\mathbf{z}_1 = \mathbf{W}_1 \mathbf{x} + \mathbf{b}_1 \quad N \times 1$$

$$\mathbf{h} = \text{ReLU}(\mathbf{z}_1) \quad N \times 1$$

$$\mathbf{z}_2 = \mathbf{W}_2 \mathbf{h} + \mathbf{b}_2 \quad V \times 1$$

$$\hat{\mathbf{y}} = \text{softmax}(\mathbf{z}_2) \quad V \times 1$$

$\mathbf{x}$

$V \times 1$

$\mathbf{h}$

$N \times 1$

$\hat{\mathbf{y}}$

$V \times 1$

# Dimensions (single input)

Column vectors

$$z_1 = W_1 x + b_1$$

$$z_1 = \begin{bmatrix} \phantom{xx} \end{bmatrix}$$
$N \times 1$

$$W_1 = \begin{bmatrix} N \times V \end{bmatrix}$$

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

$$b_1 = \begin{bmatrix} \phantom{xx} \end{bmatrix}$$
$N \times 1$

# Dimensions (single input)

Column vectors

$$z_1 = W_1 x + b_1$$

$$z_1 = \begin{bmatrix} \\ \end{bmatrix}$$
$N \times 1$

$$W_1 = \begin{bmatrix} & N \times V & \end{bmatrix}$$

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

$$b_1 = \begin{bmatrix} \\ \end{bmatrix}$$
$N \times 1$

# Dimensions (single input)

## Column vectors

$$z_1 = W_1 x + b_1$$

$z_1 = \begin{bmatrix} \phantom{xx} \end{bmatrix}$
$\quad N \times 1$

$W_1 = \begin{bmatrix} N \times V \end{bmatrix}$

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

$b_1 = \begin{bmatrix} \phantom{x} \end{bmatrix}$
$\quad N \times 1$

## Row vectors

$$z_1 = x W_1^{\top} + b_1$$

$b_1 = \begin{bmatrix} 1 \times N \end{bmatrix}$

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

$W_1 = \begin{bmatrix} N \times V \end{bmatrix}$

$b_1 = \begin{bmatrix} 1 \times N \end{bmatrix}$

# Dimensions (batch input)



Input layer     Hidden layer     Output layer

$W_1$

$B_1$

ReLU

$W_2$

$B_2$

softmax

$X$        $H$        $\hat{Y}$

# Dimensions (batch input)

Context words
vectors

$\begin{bmatrix} x^{(1)} \end{bmatrix} \cdots \begin{bmatrix} x^{(m)} \end{bmatrix}$  $\updownarrow V$

$\longleftrightarrow m$



Input layer

$W_1$

$B_1$

ReLU

**X**

Hidden layer

$W_2$

$B_2$

softmax

**H**

Output layer

**Ŷ**

# Dimensions (batch input)



Context words vectors → matrix

$$X = \left( \begin{array}{ccc} x^{(1)} & \cdots & x^{(m)} \end{array} \right) \updownarrow V$$

$\longleftrightarrow$ m

Input layer

$W_1$

$B_1$

ReLU

X

$V \times m$

Hidden layer

$W_2$

$B_2$

softmax

H

Output layer

$\hat{Y}$

# Dimensions (batch input)

Context words vectors → matrix

$$X = \begin{bmatrix} x^{(1)} & \cdots & x^{(m)} \end{bmatrix} \Big\updownarrow V$$

$\longleftrightarrow m$

**Input layer**

**Hidden layer**

**Output layer**



$W_1$    N x V

$B_1$    N x m

ReLU

$Z_1 = W_1 X + B_1$    N x m

$H = ReLU(Z_1)$    N x m

$W_2$

$B_2$

softmax

**X**

V x m

**H**

N x m

**Ŷ**

deeplearning.ai

# Dimensions (batch input)

Context words
vectors
→ matrix

$$X = \left[ \; \left| x^{(1)} \right| \cdots \left| x^{(m)} \right| \; \right] \updownarrow V$$

$\longleftrightarrow$
$m$



**Input layer**

X
V x m

**Hidden layer**

H
N x m

**Output layer**

$\hat{Y}$

$W_1$    N x V

$B_1$    N x m

ReLU

$W_2$

$B_2$

softmax

$Z_1 = W_1 X + B_1$    N x m

$H = ReLU(Z_1)$    N x m

# Dimensions (batch input)

Context words vectors
$\rightarrow$ matrix

$$X = \left( \begin{array}{ccc} \left| x^{(1)} \right| & \cdots & \left| x^{(m)} \right| \end{array} \right) \updownarrow V$$

$\longleftrightarrow m$



**Input layer**

**Hidden layer**

**Output layer**

| $W_1$ | $N \times V$ |
| $B_1$ | $N \times m$ |
| ReLU | |

| $W_2$ | |
| $B_2$ | |
| softmax | |

$Z_1 = W_1 X + B_1 \quad N \times m$

$H = \text{ReLU}(Z_1) \quad N \times m$

**X**
$V \times m$

**H**
$N \times m$

**Ŷ**
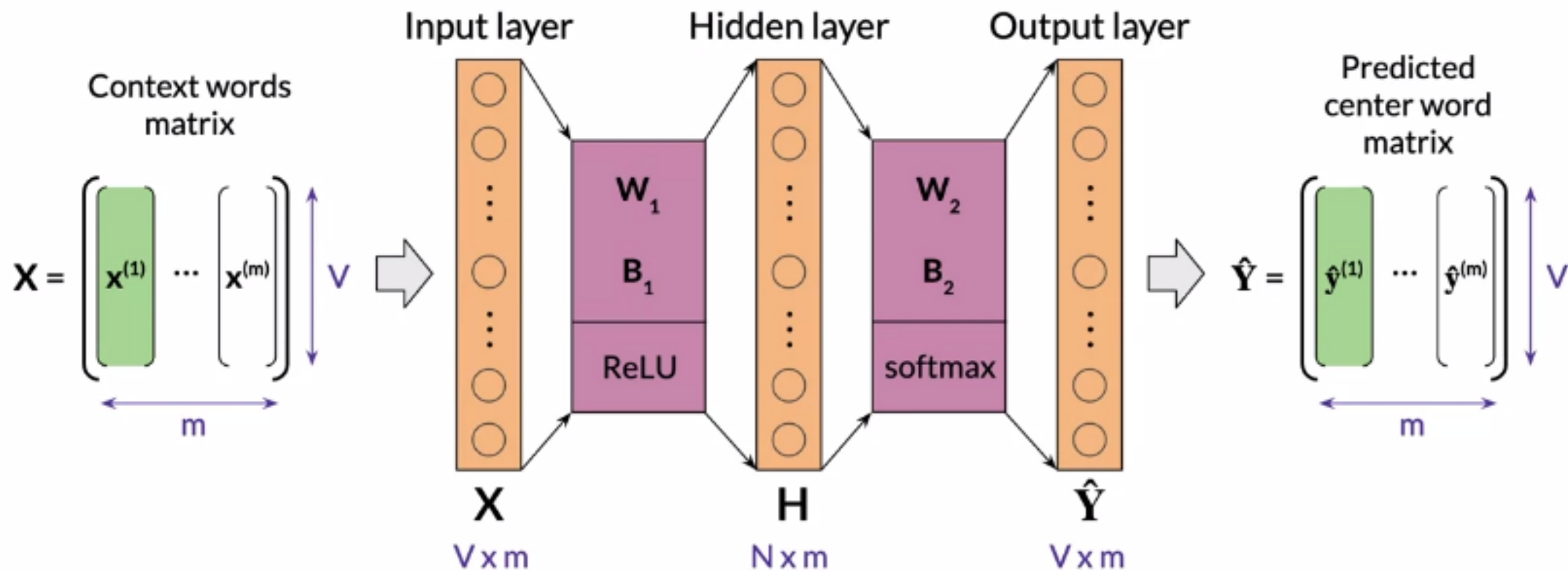
# Dimensions (batch input)

$$\begin{bmatrix} b_1 \end{bmatrix} \rightarrow B_1 = \begin{bmatrix} \begin{bmatrix} b_1 \end{bmatrix} \cdots \begin{bmatrix} b_1 \end{bmatrix} \end{bmatrix} \updownarrow N$$
$$\underset{m}{\longleftrightarrow}$$

**Input layer**　　　　**Hidden layer**　　　　**Output layer**

Context words
vectors
→ matrix

$$X = \begin{bmatrix} \begin{bmatrix} x^{(1)} \end{bmatrix} \cdots \begin{bmatrix} x^{(m)} \end{bmatrix} \end{bmatrix} \updownarrow V$$
$$\underset{m}{\longleftrightarrow}$$

$W_1$ 　 $N \times V$

$B_1$ 　 $N \times m$

ReLU

$W_2$

$B_2$

softmax

$Z_1 = W_1 X + B_1$ 　 $N \times m$

$H = \text{ReLU}(Z_1)$ 　 $N \times m$

**X**

$V \times m$

**H**

$N \times m$

**Ŷ**

# Dimensions (batch input)

$$\begin{bmatrix} b_1 \end{bmatrix} \rightarrow B_1 = \begin{bmatrix} \begin{bmatrix} b_1 \end{bmatrix} \cdots \begin{bmatrix} b_1 \end{bmatrix} \end{bmatrix} \updownarrow N \quad \textit{broadcasting}$$

Context words vectors → matrix

$$X = \begin{bmatrix} \begin{bmatrix} x^{(1)} \end{bmatrix} \cdots \begin{bmatrix} x^{(m)} \end{bmatrix} \end{bmatrix} \updownarrow V$$

Input layer

Hidden layer

Output layer

$W_1$    N x V

$B_1$    N x m

ReLU

$W_2$

$B_2$

softmax

$Z_1 = W_1 X + B_1$   N x m

$H = \text{ReLU}(Z_1)$   N x m

**X**
V x m

**H**
N x m

$\hat{Y}$

# Dimensions (batch input)

$$\begin{bmatrix} b_1 \end{bmatrix} \rightarrow \quad B_1 = \left( \begin{bmatrix} b_1 \end{bmatrix} \dots \begin{bmatrix} b_1 \end{bmatrix} \right) \updownarrow N \quad \textit{broadcasting}$$

$$\underset{m}{\longleftrightarrow}$$

**Input layer**   **Hidden layer**   **Output layer**

Context words
vectors
→ matrix

$$X = \left( \begin{bmatrix} x^{(1)} \end{bmatrix} \dots \begin{bmatrix} x^{(m)} \end{bmatrix} \right) \updownarrow V$$

$$\underset{m}{\longleftrightarrow}$$

$W_1 \qquad N \times V$

$B_1 \qquad N \times m$

ReLU

$W_2 \qquad V \times N$

$B_2 \qquad V \times m$

softmax

$Z_1 = W_1 X + B_1 \qquad N \times m$

$H = \text{ReLU}(Z_1) \qquad N \times m$

$Z_2 = W_2 H + B_2 \qquad V \times m$

$\hat{Y} = \text{softmax}(Z_2) \qquad V \times m$

**X**
$V \times m$

**H**
$N \times m$

$\hat{Y}$
$V \times m$

deeplearning.ai

# Dimensions (batch input)



Context words matrix

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}^{(1)} & \cdots & \mathbf{x}^{(m)} \end{bmatrix} \updownarrow V$$

$\xleftrightarrow{} m$

Input layer

Hidden layer

Output layer

$W_1$

$B_1$

ReLU

$W_2$

$B_2$

softmax

$\mathbf{X}$

$V \times m$

$\mathbf{H}$

$N \times m$

$\hat{\mathbf{Y}}$

$V \times m$

Predicted center word matrix

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

$\xleftrightarrow{} m$

# Dimensions (batch input)



**Context words matrix**

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}^{(1)} & \cdots & \mathbf{x}^{(m)} \end{bmatrix} \quad V$$

$m$

Input layer

Hidden layer

Output layer

$W_1$

$B_1$

ReLU

$W_2$

$B_2$

softmax

$\mathbf{X}$

$V \times m$

$\mathbf{H}$

$N \times m$

$\hat{\mathbf{Y}}$

$V \times m$

**Predicted center word matrix**

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

$m$

deeplearning.ai

# Rectified Linear Unit (ReLU)

# Rectified Linear Unit (ReLU)

Input layer

Hidden layer



$$z_1 = W_1 x + b_1$$

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

# Rectified Linear Unit (ReLU)



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

$$z_1 = W_1 x + b_1$$

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

Input layer

Hidden layer

$W_1$

$b_1$

ReLU

$x$

$h$

# Rectified Linear Unit (ReLU)

Input layer        Hidden layer



$$ReLU(x) = \max(0, x)$$

$$z_1 = W_1 x + b_1$$

$$h = ReLU(z_1)$$

# Rectified Linear Unit (ReLU)

Input layer    Hidden layer

$$z_1 = W_1 x + b_1$$

$$h = ReLU(z_1)$$

$$ReLU(x) = max(0, x)$$

# Softmax

Hidden layer          Output layer

# Softmax

Hidden layer      Output layer



$$z = W_2 h + b_2$$

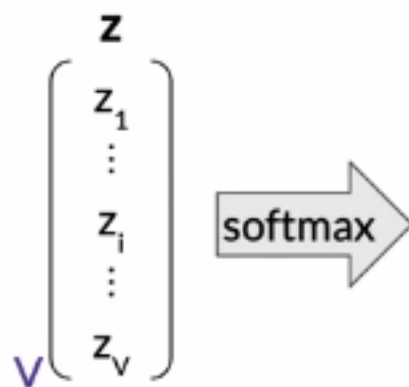$W_2$

$b_2$

softmax
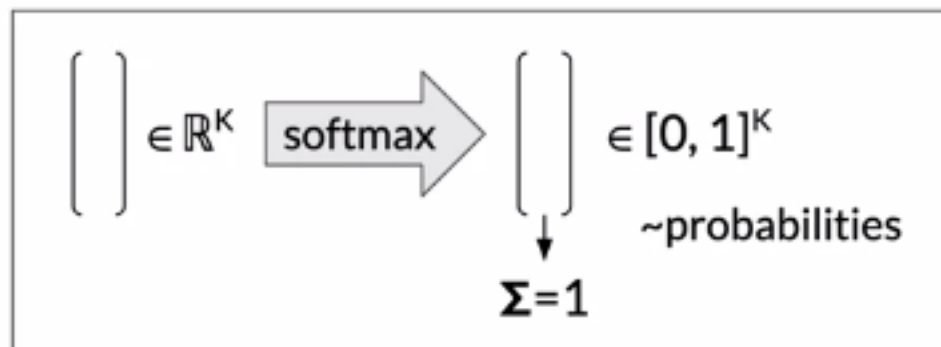
$h$      $\hat{y}$

# Softmax

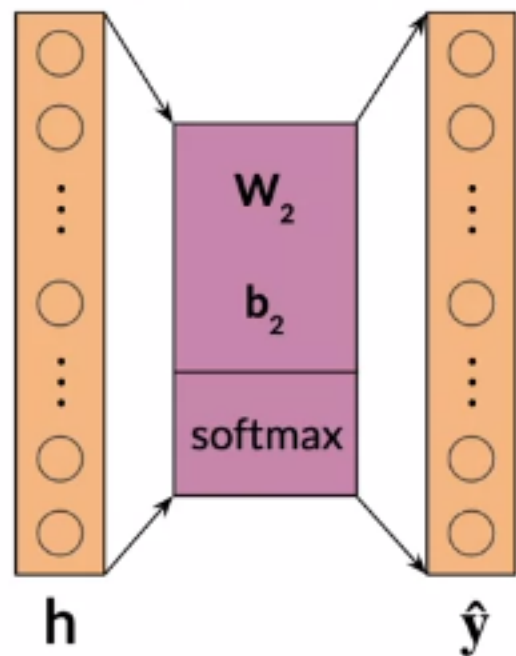Hidden layer    Output layer



$$z = W_2 h + b_2$$
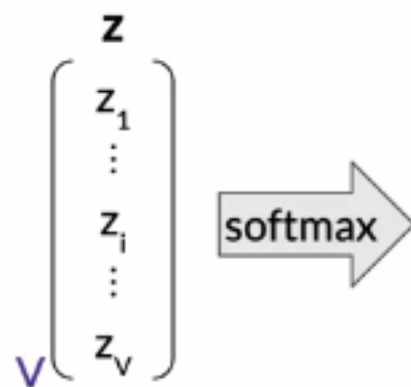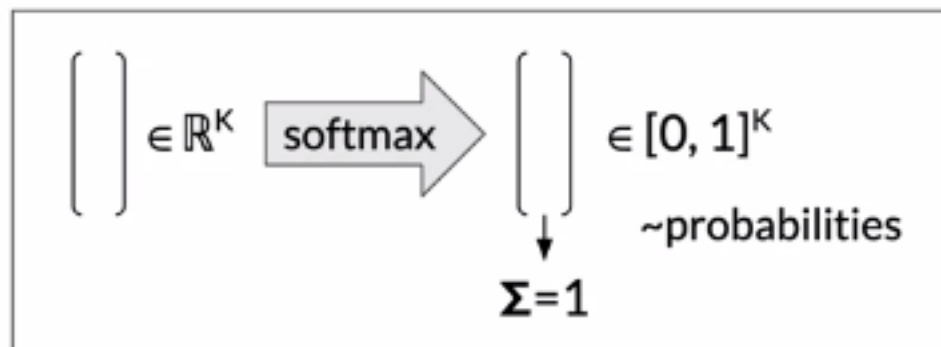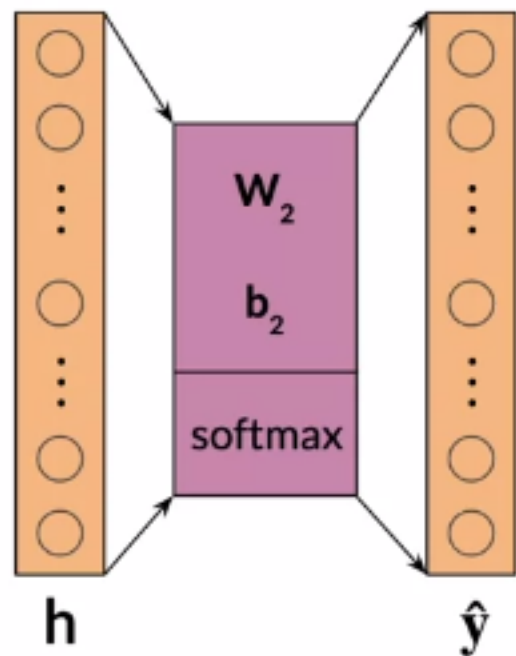
$$\hat{y} = \text{softmax}(z)$$

# Softmax



Hidden layer    Output layer

$$z = W_2 h + b_2$$

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

h    $\hat{y}$

$$\begin{bmatrix} \\ \\ \end{bmatrix} \in \mathbb{R}^K \quad \boxed{\text{softmax}} \Rightarrow \begin{bmatrix} \\ \\ \end{bmatrix} \in [0, 1]^K$$

# Softmax



Hidden layer     Output layer

$$z = W_2 h + b_2$$

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

$$\left[ \; \right] \in \mathbb{R}^K \quad \boxed{\text{softmax}} \Rightarrow \quad \left[ \; \right] \in [0, 1]^K$$

$$\Sigma = 1$$

# Softmax

Hidden layer     Output layer



$z = W_2h + b_2$

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

**h**        **ŷ**

$$\begin{bmatrix} \\ \\ \end{bmatrix} \in \mathbb{R}^K \quad \xrightarrow{\text{softmax}} \quad \begin{bmatrix} \\ \\ \end{bmatrix} \in [0, 1]^K$$

~probabilities

$\Sigma = 1$

# Softmax



**Hidden layer**   **Output layer**

$$z = W_2 h + b_2$$

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

$h$   $\hat{y}$

# Softmax

## Hidden layer    Output layer
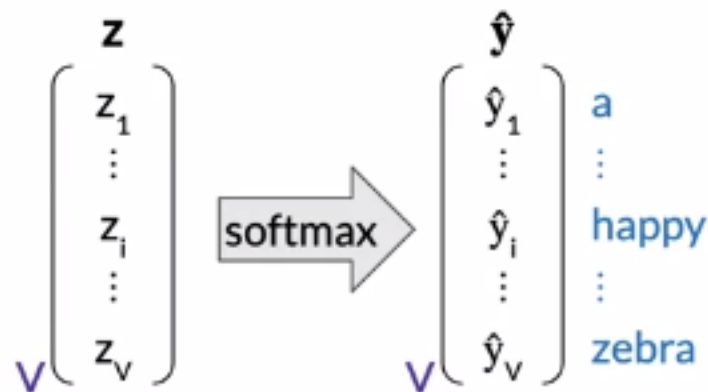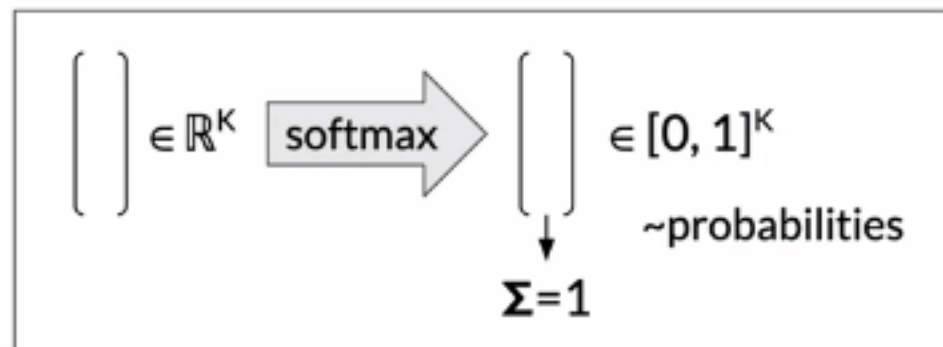


$$z = W_2 h + b_2$$

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

$$[\;] \in \mathbb{R}^K \;\xrightarrow{\text{softmax}}\; [\;] \in [0, 1]^K$$

~probabilities

$$\Sigma = 1$$

$$z \begin{bmatrix} z_1 \\ \vdots \\ z_i \\ \vdots \\ z_v \end{bmatrix} \xrightarrow{\text{softmax}}$$

# Softmax

Hidden layer     Output layer

$$z = W_2 h + b_2$$

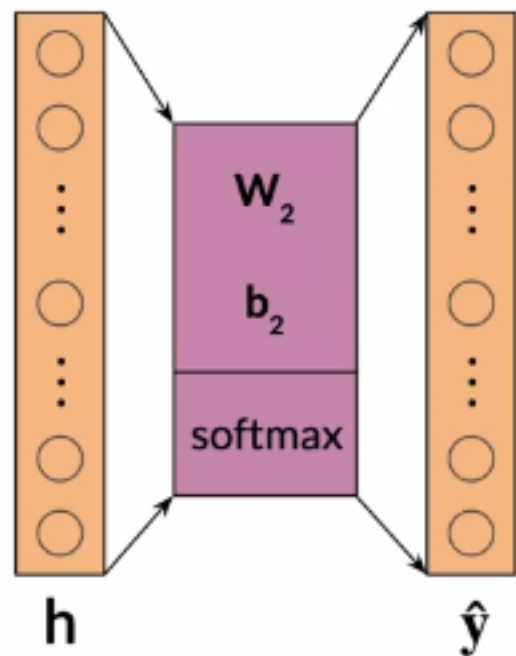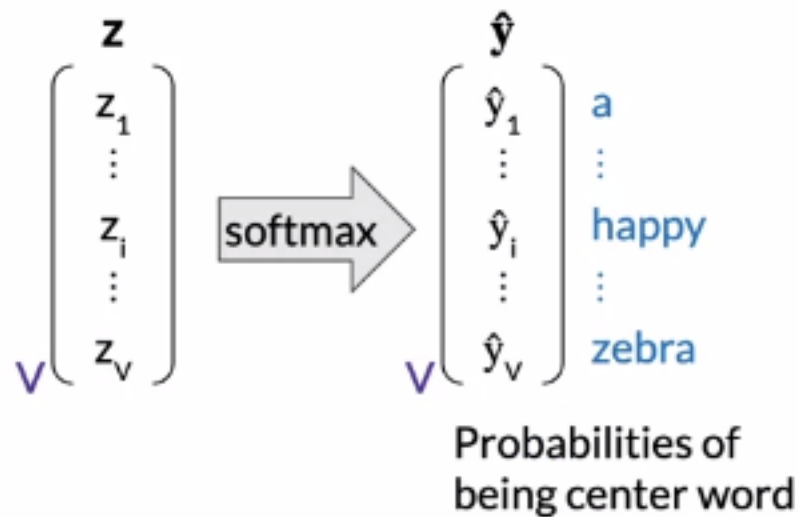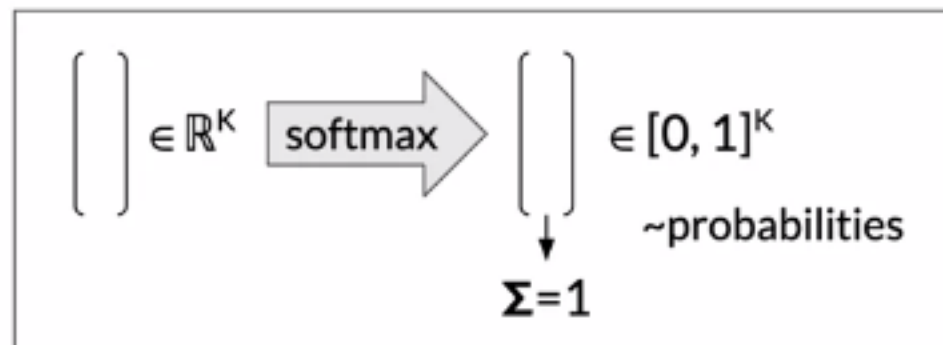$$\hat{y} = \text{softmax}(z)$$

$$\in \mathbb{R}^K \quad \boxed{\text{softmax}} \Rightarrow \quad \in [0, 1]^K$$

$$\sim \text{probabilities}$$

$$\Sigma = 1$$

$$z \quad \boxed{\text{softmax}} \Rightarrow \quad \hat{y}$$

$$\begin{array}{c} z_1 \\ \vdots \\ z_i \\ \vdots \\ z_v \end{array} \qquad \begin{array}{c} \hat{y}_1 \\ \vdots \\ \hat{y}_i \\ \vdots \\ \hat{y}_v \end{array} \qquad \begin{array}{l} a \\ \vdots \\ \text{happy} \\ \vdots \\ \text{zebra} \end{array}$$

# Softmax

**Hidden layer**   **Output layer**

$z = W_2h + b_2$

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

$$\begin{bmatrix} \ \\ \ \\ \ \end{bmatrix} \in \mathbb{R}^K \xrightarrow{\text{softmax}} \begin{bmatrix} \ \\ \ \\ \ \end{bmatrix} \in [0, 1]^K \sim \text{probabilities}$$

$\Sigma = 1$

$$\mathbf{z} \begin{bmatrix} z_1 \\ \vdots \\ z_i \\ \vdots \\ z_v \end{bmatrix} \xrightarrow{\text{softmax}} \hat{\mathbf{y}} \begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_i \\ \vdots \\ \hat{y}_v \end{bmatrix} \begin{matrix} a \\ \vdots \\ \text{happy} \\ \vdots \\ \text{zebra} \end{matrix}$$

Probabilities of
being center word

deeplearning.ai

# Softmax

**Hidden layer**      **Output layer**

$$z = W_2 h + b_2$$

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

$\mathbf{h}$      $\hat{\mathbf{y}}$

$$[\ ] \in \mathbb{R}^K \xrightarrow{\text{softmax}} [\ ] \in [0, 1]^K$$

~probabilities

$$\Sigma = 1$$

$$\mathbf{z} \begin{bmatrix} z_1 \\ \vdots \\ z_i \\ \vdots \\ z_v \end{bmatrix} \xrightarrow{\text{softmax}} \hat{\mathbf{y}} \begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_i \\ \vdots \\ \hat{y}_v \end{bmatrix} \begin{matrix} \text{a} \\ \vdots \\ \text{happy} \\ \vdots \\ \text{zebra} \end{matrix}$$

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

Probabilities of
being center word

deeplearning.ai

# Softmax

## Hidden layer     Output layer



$$z = W_2 h + b_2$$

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

$$\begin{bmatrix} \\ \\ \end{bmatrix} \in \mathbb{R}^K \quad \boxed{\text{softmax}} \Rightarrow \begin{bmatrix} \\ \\ \end{bmatrix} \in [0, 1]^K$$

~probabilities

$$\Sigma = 1$$

$$\mathbf{z} \qquad \hat{\mathbf{y}}$$

$$V\begin{bmatrix} z_1 \\ \vdots \\ z_i \\ \vdots \\ z_V \end{bmatrix} \xrightarrow{\text{softmax}} V\begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_i \\ \vdots \\ \hat{y}_V \end{bmatrix} \begin{matrix} \text{a} \\ \vdots \\ \text{happy} \\ \vdots \\ \text{zebra} \end{matrix}$$

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

Probabilities of
being center word

# Softmax: example

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

**z**

$$\begin{bmatrix} 9 \\ 8 \\ 11 \\ 10 \\ 8.5 \end{bmatrix}$$

# Softmax: example

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

| z | | exp(z) |
|---|---|---|
| 9 | | 8103 |
| 8 | | 2981 |
| 11 | exp | 59874 |
| 10 | | 22026 |
| 8.5 | | 4915 |

# Softmax: example

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

z

$$\begin{pmatrix} 9 \\ 8 \\ 11 \\ 10 \\ 8.5 \end{pmatrix}$$

exp →

exp(z)

$$\begin{pmatrix} 8103 \\ 2981 \\ 59874 \\ 22026 \\ 4915 \end{pmatrix}$$

Σ=97899

# Softmax: example

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

| z | | exp(z) |
|---|---|---|
| 9 | | 8103 |
| 8 | | 2981 |
| 11 | exp | 59874 |
| 10 | | 22026 |
| 8.5 | | 4915 |

Σ=97899

# Softmax: example

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

**z**

$$\begin{pmatrix} 9 \\ 8 \\ 11 \\ 10 \\ 8.5 \end{pmatrix}$$

→ exp →

**exp(z)**

$$\begin{pmatrix} 8103 \\ 2981 \\ 59874 \\ 22026 \\ 4915 \end{pmatrix}$$

→ /Σ →

**ŷ = softmax(z)**

$$\begin{pmatrix} \\ \\ \\ \\ \end{pmatrix}$$

Σ=97899

# Softmax: example

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

| z | | exp(z) | | $\hat{y}$ = softmax(z) |
|---|---|---|---|---|
| 9 | | 8103 | | 0.083 |
| 8 | exp | 2981 | /Σ | |
| 11 | | 59874 | | |
| 10 | | 22026 | | |
| 8.5 | | 4915 | | |

Σ=97899

deeplearning.ai

# Softmax: example

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

| z | exp(z) | $\hat{y}$ = softmax(z) |
|---|--------|------------------------|
| 9 | 8103 | 0.083 |
| 8 | 2981 | |
| 11 | 59874 | |
| 10 | 22026 | |
| 8.5 | 4915 | |

exp →   / Σ →

Σ=97899

# Softmax: example

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

| z | | exp(z) | | ŷ = softmax(z) |
|---|---|---|---|---|
| 9 | | 8103 | | 0.083 |
| 8 | | 2981 | | 0.03 |
| 11 | exp → | 59874 | /Σ → | 0.612 |
| 10 | | 22026 | | 0.225 |
| 8.5 | | 4915 | | 0.05 |

Σ=97899

# Softmax: example

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

| z | | exp(z) | | ŷ = softmax(z) |
|---|---|---|---|---|
| 9 | | 8103 | | 0.083 |
| 8 | exp | 2981 | /Σ | 0.03 |
| 11 | | 59874 | | 0.612 |
| 10 | | 22026 | | 0.225 |
| 8.5 | | 4915 | | 0.05 |

Σ=97899

Σ=1

deeplearning.ai

# Softmax: example

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

| z | | exp(z) | | $\hat{y}$ = softmax(z) | |
|---|---|---|---|---|---|
| 9 | | 8103 | | 0.083 | am |
| 8 | exp | 2981 | /Σ | 0.03 | because |
| 11 | | 59874 | | 0.612 | happy |
| 10 | | 22026 | | 0.225 | I |
| 8.5 | | 4915 | | 0.05 | learning |

Σ=97899

Σ=1

# Softmax: example

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

| z | | exp(z) | | ŷ = softmax(z) | |
|---|---|---|---|---|---|
| 9 | | 8103 | | 0.083 | am |
| 8 | exp | 2981 | /Σ | 0.03 | because |
| 11 | | 59874 | | 0.612 | happy |
| 10 | | 22026 | | 0.225 | I |
| 8.5 | | 4915 | | 0.05 | learning |

Σ=97899

Σ=1

# Softmax: example

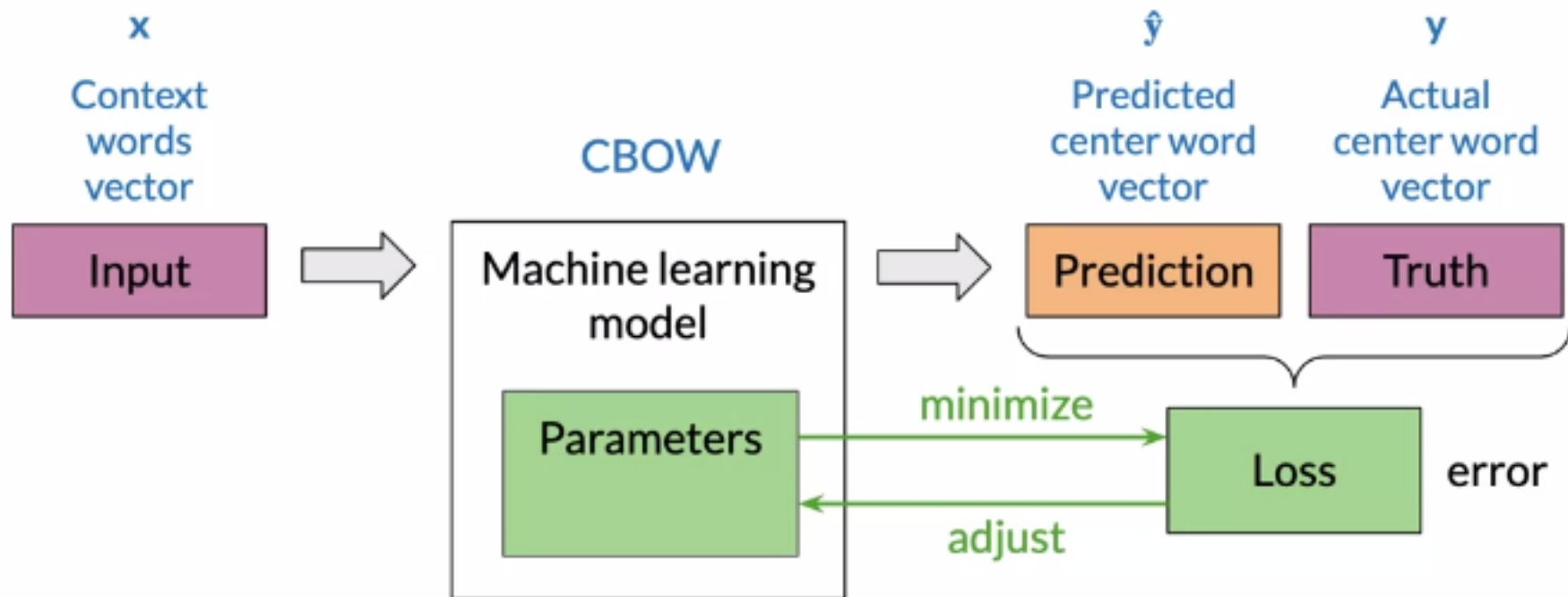$$\hat{y}_i = \frac{e^{z_i}}{\sum_{j=1}^{V} e^{z_j}}$$

| z | | exp(z) | | ŷ = softmax(z) | |
|---|---|---|---|---|---|
| 9 | | 8103 | | 0.083 | am |
| 8 | exp | 2981 | /Σ | 0.03 | because |
| 11 | | 59874 | | 0.612 | happy |
| 10 | | 22026 | | 0.225 | I |
| 8.5 | | 4915 | | 0.05 | learning |

Σ=97899

Σ=1

# Softmax: example

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

| z | | exp(z) | | $\hat{y}$ = softmax(z) | |
|---|---|---|---|---|---|
| 9 | | 8103 | | 0.083 | am |
| 8 | | 2981 | | 0.03 | because |
| 11 | exp | 59874 | /Σ | 0.612 | happy ← Predicted center word |
| 10 | | 22026 | | 0.225 | I |
| 8.5 | | 4915 | | 0.05 | learning |

Σ=97899

Σ=1

# Loss

# Loss

# Loss

# Loss

# Loss

# Loss

# Loss

# Cross-entropy loss

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

deeplearning.ai

# 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
$$\hat{\mathbf{y}} = \begin{pmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_V \end{pmatrix}$$

# 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

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

**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

**y**

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

am
because
happy
I
learning

# 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

| **y** | | **ŷ** |
|---|---|---|
| 0 | am | 0.083 |
| 0 | because | 0.03 |
| 1 | happy | 0.611 |
| 0 | I | 0.225 |
| 0 | learning | 0.05 |

# 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 $\hat{\mathbf{y}} = \begin{pmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_V \end{pmatrix}$
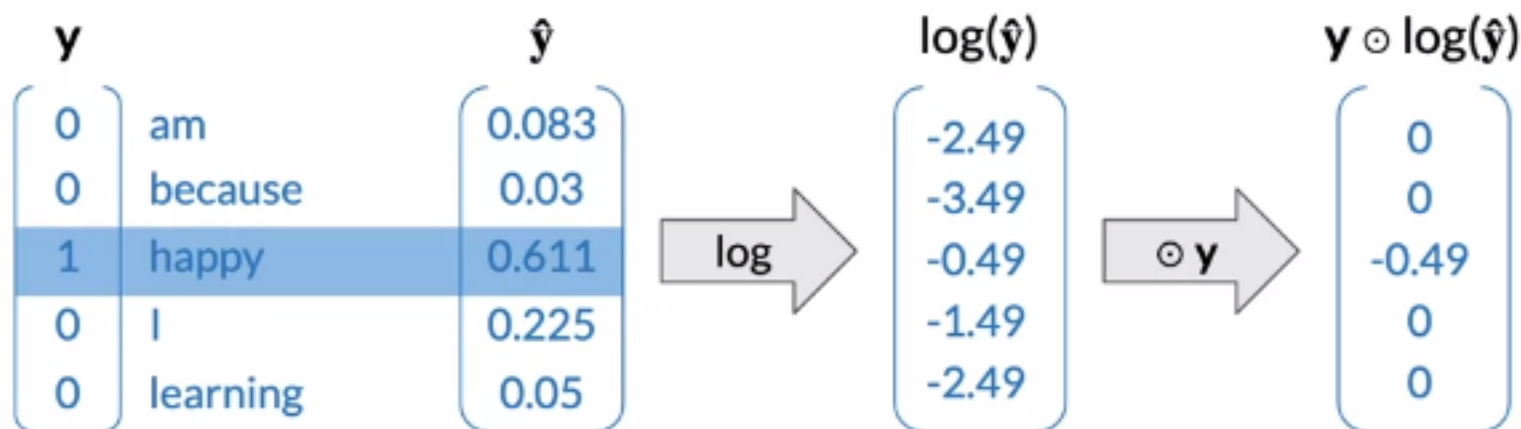
I am happy because I am learning

| $\mathbf{y}$ | | $\hat{\mathbf{y}}$ |
|---|---|---|
| 0 | am | 0.083 |
| 0 | because | 0.03 |
| 1 | happy | 0.611 |
| 0 | I | 0.225 |
| 0 | learning | 0.05 |

# 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

| **y** | | **ŷ** | | **log(ŷ)** |
|---|---|---|---|---|
| 0 | am | 0.083 | | -2.49 |
| 0 | because | 0.03 | | -3.49 |
| 1 | happy | 0.611 | log | -0.49 |
| 0 | I | 0.225 | | -1.49 |
| 0 | learning | 0.05 | | -2.49 |

# Cross-entropy loss

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

Actual $\quad \mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_V \end{bmatrix}$
 $\qquad$ Predicted $\quad \hat{\mathbf{y}} = \begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_V \end{bmatrix}$
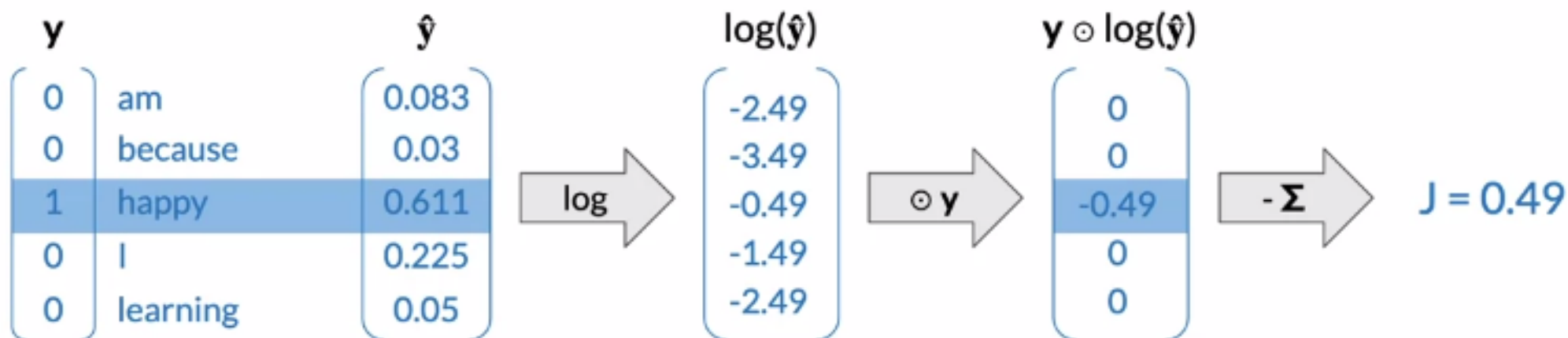
I am happy because I am learning

| **y** | | **ŷ** | | **log(ŷ)** | | **y ⊙ log(ŷ)** |
|---|---|---|---|---|---|---|
| 0 | am | 0.083 | | -2.49 | | 0 |
| 0 | because | 0.03 | | -3.49 | | 0 |
| 1 | happy | 0.611 | log → | -0.49 | ⊙ **y** → | -0.49 |
| 0 | I | 0.225 | | -1.49 | | 0 |
| 0 | learning | 0.05 | | -2.49 | | 0 |

deeplearning.ai

# 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

| **y** | | **ŷ** | | log(ŷ) | | **y** ⊙ log(ŷ) |
|---|---|---|---|---|---|---|
| 0 | am | 0.083 | | -2.49 | | 0 |
| 0 | because | 0.03 | | -3.49 | | 0 |
| 1 | happy | 0.611 | log → | -0.49 | ⊙ y → | -0.49 |
| 0 | I | 0.225 | | -1.49 | | 0 |
| 0 | learning | 0.05 | | -2.49 | | 0 |

# 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

| $\mathbf{y}$ | | $\hat{\mathbf{y}}$ | | $\log(\hat{\mathbf{y}})$ | | $\mathbf{y} \odot \log(\hat{\mathbf{y}})$ | |
|---|---|---|---|---|---|---|---|
| 0 | am | 0.083 | | -2.49 | | 0 | |
| 0 | because | 0.03 | | -3.49 | | 0 | |
| 1 | happy | 0.611 | log | -0.49 | $\odot \mathbf{y}$ | -0.49 | $-\Sigma$ |
| 0 | I | 0.225 | | -1.49 | | 0 | |
| 0 | learning | 0.05 | | -2.49 | | 0 | |

$J = 0.49$

# Cross-entropy loss

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

**y**

$$\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$$

am
because
happy
I
learning

# Cross-entropy loss

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

| **y** | | **ŷ** |
|---|---|---|
| 0 | am | 0.96 |
| 0 | because | 0.01 |
| 1 | happy | 0.01 |
| 0 | I | 0.01 |
| 0 | learning | 0.01 |

# Cross-entropy loss

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

| **y** | | **ŷ** | log(ŷ) |
|---|---|---|---|
| 0 | am | 0.96 | -0.04 |
| 0 | because | 0.01 | -4.61 |
| 1 | happy | 0.01 | -4.61 |
| 0 | I | 0.01 | -4.61 |
| 0 | learning | 0.01 | -4.61 |

log

# Cross-entropy loss

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

| **y** | | **ŷ** |
|---|---|---|
| 0 | am | 0.96 |
| 0 | because | 0.01 |
| 1 | happy | 0.01 |
| 0 | I | 0.01 |
| 0 | learning | 0.01 |

**log →**

**log(ŷ)**

| |
|---|
| -0.04 |
| -4.61 |
| -4.61 |
| -4.61 |
| -4.61 |

**⊙ y →**

**y ⊙ log(ŷ)**

| |
|---|
| 0 |
| 0 |
| -4.61 |
| 0 |
| 0 |

deeplearning.ai

# Cross-entropy loss

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

**y**

$$\begin{bmatrix} 0 & \text{am} \\ 0 & \text{because} \\ 1 & \text{happy} \\ 0 & \text{I} \\ 0 & \text{learning} \end{bmatrix}$$

**ŷ**

$$\begin{bmatrix} 0.96 \\ 0.01 \\ 0.01 \\ 0.01 \\ 0.01 \end{bmatrix}$$

log →

**log(ŷ)**

$$\begin{bmatrix} -0.04 \\ -4.61 \\ -4.61 \\ -4.61 \\ -4.61 \end{bmatrix}$$

⊙ **y** →

**y ⊙ log(ŷ)**

$$\begin{bmatrix} 0 \\ 0 \\ -4.61 \\ 0 \\ 0 \end{bmatrix}$$

- Σ →

J = 4.61

# Cross-entropy loss

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

| **y** | | **ŷ** |
|:---:|:---:|:---:|
| 0 | am | 0.96 |
| 0 | because | 0.01 |
| 1 | happy | 0.01 |
| 0 | I | 0.01 |
| 0 | learning | 0.01 |

log →

$\log(\hat{y})$

-0.04
-4.61
-4.61
-4.61
-4.61

⊙ **y** →

$y \odot \log(\hat{y})$

0
0
-4.61
0
0

-Σ →

J = 4.61

>

J(correct) = 0.49

# Cross-entropy loss

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

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

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

$\rightarrow J = 4.61$

deeplearning.ai

# Cross-entropy loss

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

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

| 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

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

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

| **y** | | **ŷ** |
|---|---|---|
| 0 | am | 0.96 |
| 0 | because | 0.01 |
| 1 | happy | 0.01 |
| 0 | I | 0.01 |
| 0 | learning | 0.01 |

→ J = 4.61

incorrect predictions: penalty

correct predictions: reward

cross-entropy loss

$\hat{y}_{\text{actual word}}$

# Training process

- Forward propagation

# Training process

- Forward propagation

- Cost

# Training process

- Forward propagation

- Cost

- Backpropagation and gradient descent

# Forward propagation

# Forward propagation

# Forward propagation

# Forward propagation

Forward propagation

$$Z_1 = W_1 X + B_1 \qquad Z_2 = W_2 H + B_2$$
$$H = \text{ReLU}(Z_1) \qquad \hat{Y} = \text{softmax}(Z_2)$$

Input layer    Hidden layer    Output layer

Context words matrix

$$X = \begin{pmatrix} \begin{pmatrix} x^{(1)} \end{pmatrix} \cdots \begin{pmatrix} x^{(m)} \end{pmatrix} \end{pmatrix}$$

$W_1$

$B_1$

ReLU

$W_2$

$B_2$

softmax

**X**      **H**      **Ŷ**

Predicted center word matrix

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

# Cost

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

# Cost

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

Cost: mean of losses

# Cost

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

Cost: mean of losses

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

Predicted center word matrix

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

Actual center word matrix

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

# Cost

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

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

$$\mathbf{\hat{Y}} = \left[ \left[ \hat{\mathbf{y}}^{(1)} \right] \cdots \left[ \hat{\mathbf{y}}^{(m)} \right] \right]$$

Actual center word matrix

$$\mathbf{Y} = \left[ \left[ \mathbf{y}^{(1)} \right] \cdots \left[ \mathbf{y}^{(m)} \right] \right]$$

# Minimizing the cost

# Minimizing the cost

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

$$\frac{\partial J_{batch}}{\partial \mathbf{W_1}}, \frac{\partial J_{batch}}{\partial \mathbf{W_2}}, \frac{\partial J_{batch}}{\partial \mathbf{b_1}}, \frac{\partial J_{batch}}{\partial \mathbf{b_2}}$$

# 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

# Backpropagation

$$\frac{\partial J_{batch}}{\partial \mathbf{W_1}}$$

$$\frac{\partial J_{batch}}{\partial \mathbf{W_2}}$$

$$\frac{\partial J_{batch}}{\partial \mathbf{b_1}}$$

$$\frac{\partial J_{batch}}{\partial \mathbf{b_2}}$$

# Backpropagation

$$\frac{\partial J_{batch}}{\partial \mathbf{W_1}} = \frac{1}{m} \text{ReLU} \left( \mathbf{W_2}^\top (\hat{\mathbf{Y}} - \mathbf{Y}) \right) \mathbf{X}^\top$$

$$\frac{\partial J_{batch}}{\partial \mathbf{W_2}}$$

$$\frac{\partial J_{batch}}{\partial \mathbf{b_1}}$$

$$\frac{\partial J_{batch}}{\partial \mathbf{b_2}}$$

# Backpropagation

$$\frac{\partial J_{batch}}{\partial \mathbf{W_1}} = \frac{1}{m} \text{ReLU} \left( \mathbf{W_2^\top} (\hat{\mathbf{Y}} - \mathbf{Y}) \right) \mathbf{X}^\top$$

$$\frac{\partial J_{batch}}{\partial \mathbf{W_2}} = \frac{1}{m} (\hat{\mathbf{Y}} - \mathbf{Y}) \mathbf{H}^\top$$

$$\frac{\partial J_{batch}}{\partial \mathbf{b_1}}$$

$$\frac{\partial J_{batch}}{\partial \mathbf{b_2}}$$

# Backpropagation

$$\frac{\partial J_{batch}}{\partial \mathbf{W_1}} = \frac{1}{m} \text{ReLU}\left(\mathbf{W_2^\top}(\hat{\mathbf{Y}} - \mathbf{Y})\right)\mathbf{X}^\top$$

$$\frac{\partial J_{batch}}{\partial \mathbf{W_2}} = \frac{1}{m}(\hat{\mathbf{Y}} - \mathbf{Y})\mathbf{H}^\top$$

$$\frac{\partial J_{batch}}{\partial \mathbf{b_1}} = \frac{1}{m} \text{ReLU}\left(\mathbf{W_2^\top}(\hat{\mathbf{Y}} - \mathbf{Y})\right)\mathbf{1}_m^\top$$

$$\frac{\partial J_{batch}}{\partial \mathbf{b_2}}$$

# Backpropagation

$$\frac{\partial J_{batch}}{\partial \mathbf{W_1}} = \frac{1}{m} \text{ReLU}\left(\mathbf{W_2^\top}(\hat{\mathbf{Y}} - \mathbf{Y})\right)\mathbf{X}^\top$$

$$\frac{\partial J_{batch}}{\partial \mathbf{W_2}} = \frac{1}{m}(\hat{\mathbf{Y}} - \mathbf{Y})\mathbf{H}^\top$$

$$\frac{\partial J_{batch}}{\partial \mathbf{b_1}} = \frac{1}{m} \text{ReLU}\left(\mathbf{W_2^\top}(\hat{\mathbf{Y}} - \mathbf{Y})\right)\mathbf{1}_m^\top$$

$$\frac{\partial J_{batch}}{\partial \mathbf{b_2}}$$

# Backpropagation

$$\frac{\partial J_{batch}}{\partial \mathbf{W_1}} = \frac{1}{m} \text{ReLU}\left(\mathbf{W_2}^\top (\hat{\mathbf{Y}} - \mathbf{Y})\right) \mathbf{X}^\top$$

$$\frac{\partial J_{batch}}{\partial \mathbf{W_2}} = \frac{1}{m} (\hat{\mathbf{Y}} - \mathbf{Y}) \mathbf{H}^\top$$

$$\frac{\partial J_{batch}}{\partial \mathbf{b_1}} = \frac{1}{m} \text{ReLU}\left(\mathbf{W_2}^\top (\hat{\mathbf{Y}} - \mathbf{Y})\right) \mathbf{1}_m^\top$$

$$\frac{\partial J_{batch}}{\partial \mathbf{b_2}}$$



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

# Backpropagation

$$\frac{\partial J_{batch}}{\partial \mathbf{W_1}} = \frac{1}{m}\text{ReLU}\left(\mathbf{W_2^\top}(\hat{\mathbf{Y}} - \mathbf{Y})\right)\mathbf{X^\top}$$

$$\frac{\partial J_{batch}}{\partial \mathbf{W_2}} = \frac{1}{m}(\hat{\mathbf{Y}} - \mathbf{Y})\mathbf{H^\top}$$

$$\frac{\partial J_{batch}}{\partial \mathbf{b_1}} = \frac{1}{m}\text{ReLU}\left(\mathbf{W_2^\top}(\hat{\mathbf{Y}} - \mathbf{Y})\right)\mathbf{1}_m^\top$$

$$\frac{\partial J_{batch}}{\partial \mathbf{b_2}} = \frac{1}{m}(\hat{\mathbf{Y}} - \mathbf{Y})\mathbf{1}_m^\top$$

$$\mathbf{1}_m = \begin{bmatrix} 1, ..., 1 \end{bmatrix}$$

$$\overset{m}{\longleftrightarrow}$$

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

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

# Gradient descent

Hyperparameter: learning rate $\alpha$

$$\mathbf{W_1} := \mathbf{W_1} - \alpha \frac{\partial J_{batch}}{\partial \mathbf{W_1}}$$

$$\mathbf{W_2} := \mathbf{W_2} - \alpha \frac{\partial J_{batch}}{\partial \mathbf{W_2}}$$

$$\mathbf{b_1} := \mathbf{b_1} - \alpha \frac{\partial J_{batch}}{\partial \mathbf{b_1}}$$

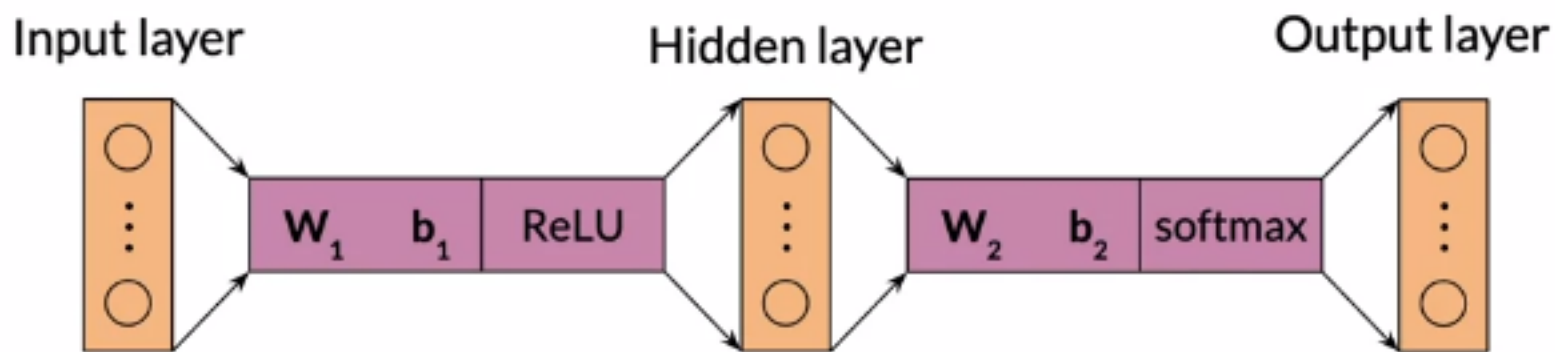$$\mathbf{b_2} := \mathbf{b_2} - \alpha \frac{\partial J_{batch}}{\partial \mathbf{b_2}}$$

# Gradient descent

Hyperparameter: learning rate α

$$\mathbf{W_1} := \mathbf{W_1} - \alpha \frac{\partial J_{batch}}{\partial \mathbf{W_1}}$$

$$\mathbf{W_2} := \mathbf{W_2} - \alpha \frac{\partial J_{batch}}{\partial \mathbf{W_2}}$$

$$\mathbf{b_1} := \mathbf{b_1} - \alpha \frac{\partial J_{batch}}{\partial \mathbf{b_1}}$$

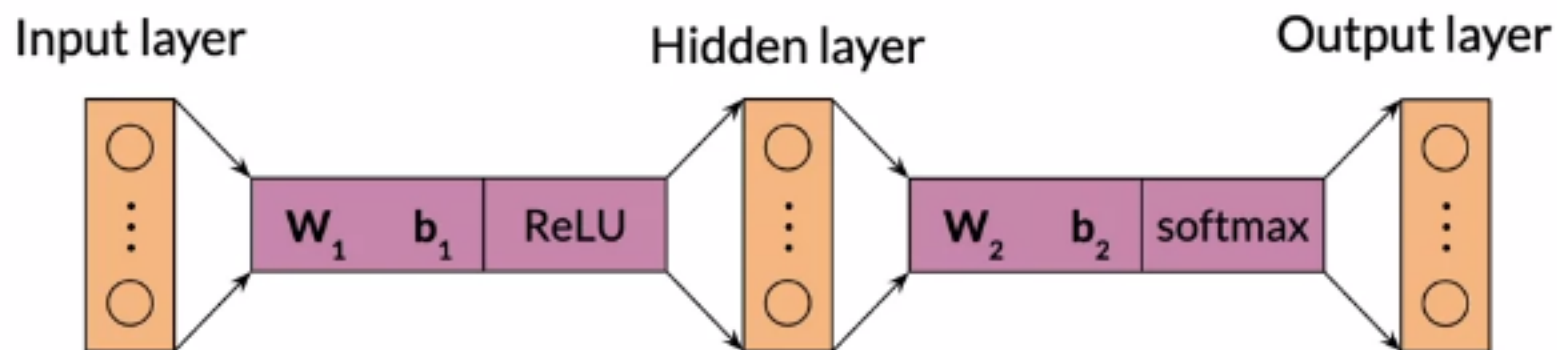$$\mathbf{b_2} := \mathbf{b_2} - \alpha \frac{\partial J_{batch}}{\partial \mathbf{b_2}}$$

# Extracting word embedding vectors: option 1

# Extracting word embedding vectors: option 1



Input layer       Hidden layer       Output layer

$W_1$   $b_1$   ReLU

$W_2$   $b_2$   softmax

$$W_1 = \left[ \left[ w^{(1)} \right] \quad \cdots \quad \left[ w^{(V)} \right] \right] \updownarrow N$$

$\longleftrightarrow V$

# Extracting word embedding vectors: option 1



Input layer        Hidden layer       Output layer

$W_1$   $b_1$   ReLU

$W_2$   $b_2$   softmax

$$W_1 = \left[ \left[ w^{(1)} \right] \; \cdots \; \left[ w^{(V)} \right] \right] \updownarrow N$$

$$\underset{V}{\longleftrightarrow}$$

$$x = \begin{bmatrix} \\ \\ \\ \\ \\ \end{bmatrix} \begin{matrix} \text{am} \\ \text{because} \\ \text{happy} \\ \text{I} \\ \text{learning} \end{matrix} \updownarrow V$$
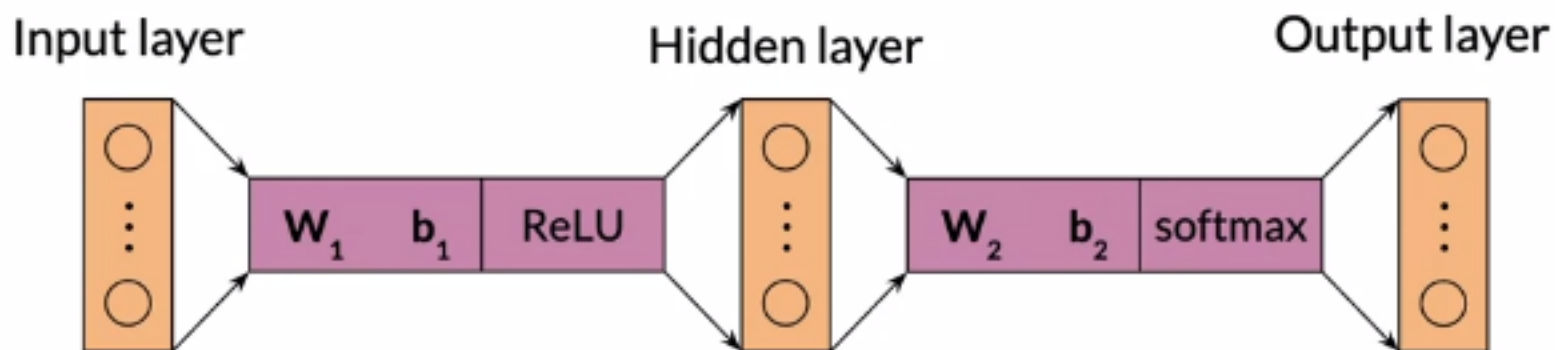
deeplearning.ai

# Extracting word embedding vectors: option 1



$$W_1 = \begin{bmatrix} \overset{\text{am}}{\boxed{w^{(1)}}} & \cdots & w^{(V)} \end{bmatrix} \updownarrow N$$

$$\underset{V}{\longleftrightarrow}$$

$$x = \begin{bmatrix} \boxed{\text{am}} \\ \text{because} \\ \text{happy} \\ \text{I} \\ \text{learning} \end{bmatrix} \updownarrow V$$

# Extracting word embedding vectors: option 2



Input layer · Hidden layer · Output layer

$W_1$ $b_1$ ReLU $W_2$ $b_2$ softmax

$$W_2 = \begin{bmatrix} \left[\begin{array}{c} w^{(1)} \end{array}\right] \\ \dots \\ \left[\begin{array}{c} w^{(V)} \end{array}\right] \end{bmatrix} \; V$$

$N$

# Extracting word embedding vectors: option 2



Input layer — $W_1$ $b_1$ — ReLU — Hidden layer — $W_2$ $b_2$ — softmax — Output layer

$$W_2 = \begin{bmatrix} \begin{bmatrix} w^{(1)} \end{bmatrix} \\ \dots \\ \begin{bmatrix} w^{(V)} \end{bmatrix} \end{bmatrix} \Bigg\} V$$

$$\underbrace{\qquad}_{N}$$

$$x = \begin{bmatrix} \text{am} \\ \text{because} \\ \text{happy} \\ \text{I} \\ \text{learning} \end{bmatrix} \Bigg\} V$$
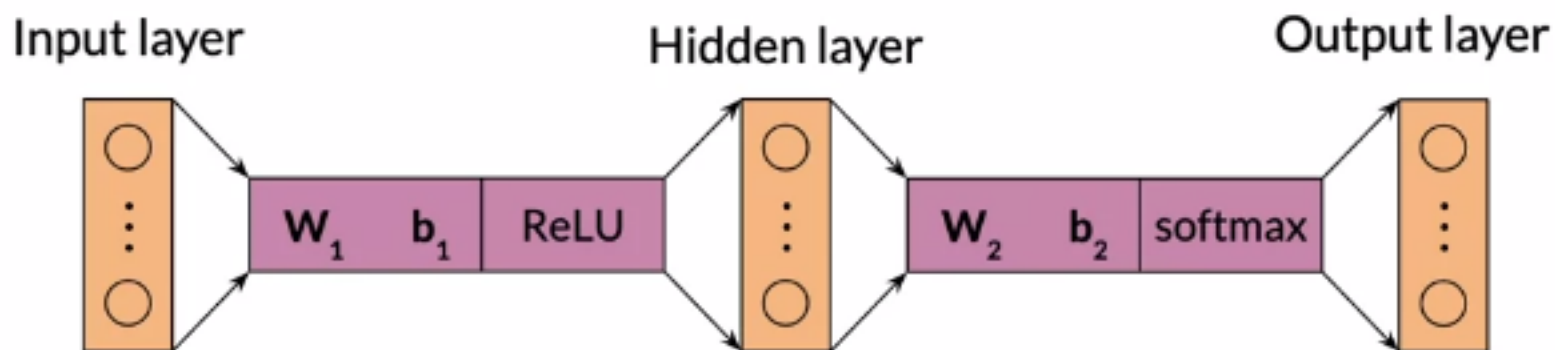
# Extracting word embedding vectors: option 3

$$W_1 = \left[ \begin{pmatrix} w_1^{(1)} \end{pmatrix} \quad \cdots \quad \begin{pmatrix} w_1^{(V)} \end{pmatrix} \right] \qquad W_2 = \begin{bmatrix} \begin{pmatrix} w_2^{(1)} \end{pmatrix} \\ \cdots \\ \begin{pmatrix} w_2^{(V)} \end{pmatrix} \end{bmatrix}$$

# Extracting word embedding vectors: option 3

$$W_1 = \left[ \left( w_1^{(1)} \right) \quad \cdots \quad \left( w_1^{(V)} \right) \right] \qquad W_2 = \left[ \begin{array}{c} \left( w_2^{(1)} \right) \\ \cdots \\ \left( w_2^{(V)} \right) \end{array} \right]$$

$$W_3 = 0.5\, (W_1 + W_2^T) = \left[ \left( w_3^{(1)} \right) \quad \cdots \quad \left( w_3^{(V)} \right) \right] \quad N$$

$$V$$

# Extracting word embedding vectors: option 3

$$W_1 = \begin{bmatrix} \boxed{w_1^{(1)}} & \cdots & \begin{bmatrix} w_1^{(V)} \end{bmatrix} \end{bmatrix} \qquad W_2 = \begin{bmatrix} \boxed{w_2^{(1)}} \\ \cdots \\ \begin{bmatrix} w_2^{(V)} \end{bmatrix} \end{bmatrix}$$

$$W_3 = 0.5\,(W_1 + W_2^T) = \begin{bmatrix} \boxed{w_3^{(1)}} & \cdots & \begin{bmatrix} w_3^{(V)} \end{bmatrix} \end{bmatrix} \Big\} N$$

$\underbrace{\qquad\qquad}_{V}$

$$x = \begin{bmatrix} \begin{array}{l} \text{am} \\ \text{because} \\ \text{happy} \\ \text{I} \\ \text{learning} \end{array} \end{bmatrix} \updownarrow V$$

# Intrinsic evaluation

# Intrinsic evaluation

Test relationships between words

# Intrinsic evaluation

Test relationships between words
- Analogies

# Intrinsic evaluation

Test relationships between words

- Analogies

   **Semantic analogies**

# Intrinsic evaluation

Test relationships between words
- Analogies

  **Semantic analogies**

  "France" is to "Paris" as "Italy" is to <?>

  **Syntactic analogies**

# 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 <?>

# 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

# 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

| 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
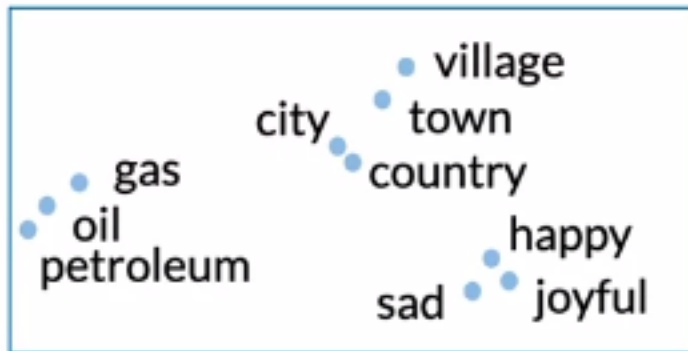
Test relationships between words

- Analogies

- Clustering

# 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



deeplearning.ai

# Intrinsic evaluation

Test relationships between words

- Analogies

- Clustering

- Visualization

# Extrinsic evaluation

Test word embeddings on external task

# Extrinsic evaluation

Test word embeddings on external task
e.g. named entity recognition, parts-of-speech tagging

# Extrinsic evaluation

Test word embeddings on external task
e.g. named entity recognition, parts-of-speech tagging

# 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*

deeplearning.ai

# Extrinsic evaluation

Test word embeddings on external task
e.g. named entity recognition, parts-of-speech tagging

# 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

# Extrinsic evaluation

Test word embeddings on external task
e.g. named entity recognition, parts-of-speech tagging

\+    Evaluates actual usefulness of embeddings

\-    Time-consuming

# 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

deeplearning.ai

Conclusion

# Recap and assignment

# Recap and assignment

- Data preparation

- Word representations

# Recap and assignment

- Data preparation

- Word representations

- Continuous bag-of-words model

# Recap and assignment

- Data preparation

- Word representations

- Continuous bag-of-words model

- Evaluation

# Going further

- Advanced language modelling and word embeddings

# Going further

- Advanced language modelling and word embeddings

- NLP and machine learning libraries

# 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)
```