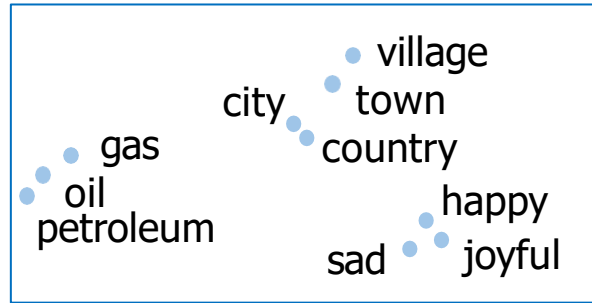# Word embedding with neural network

- Identify the key concepts of word representations

- Generate word embeddings

- Prepare text for machine learning

- Implement the continuous bag-of-words model

# Some basic applications of word embeddings

Semantic analogies and similarity
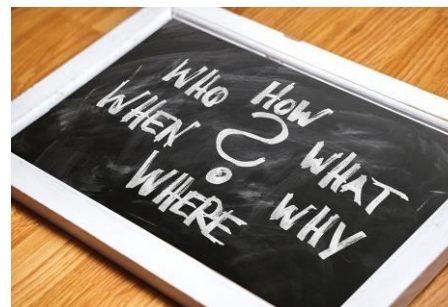
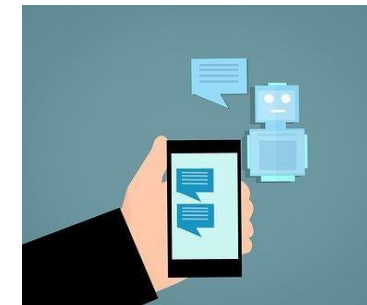Sentiment analysis

Classification of customer feedback

## Advanced applications of word embeddings

Machine translation

Information extraction

Question answering
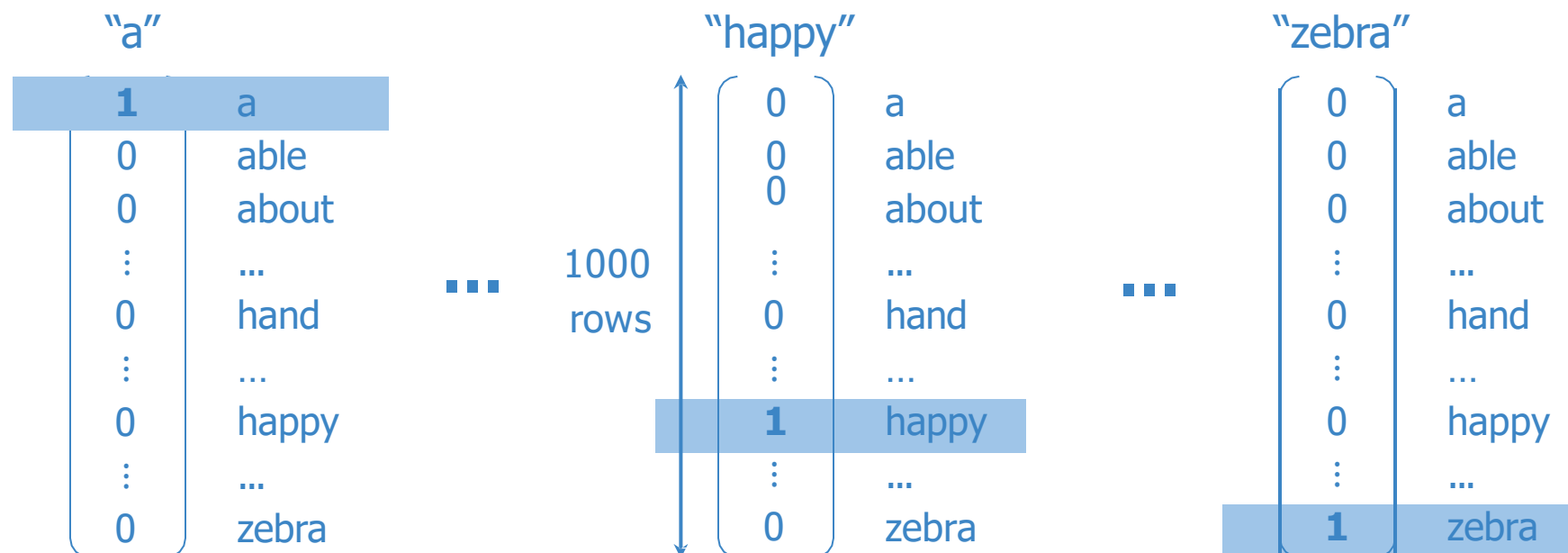
# Basic Word  Representations

o   Integers

o   One-hot vectors

o   Word embeddings

- Integers

  +  Simple

  -   Ordering: little semantic sense

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

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

1000

615 **?!** 621 **?!**

# 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 |

- + No implied ordering

- + Simple

- - Huge vectors

-  - No embedded meaning

# 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 |

# Meaning as vectors

# Word embedding vectors

"happy"

$+$ Low dimension

~100—~1000
rows

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

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

  forest ≈ tree     forest ≉ ticket

  ○ e.g. analogies

  Paris:France :: Rome:?

# Terminology

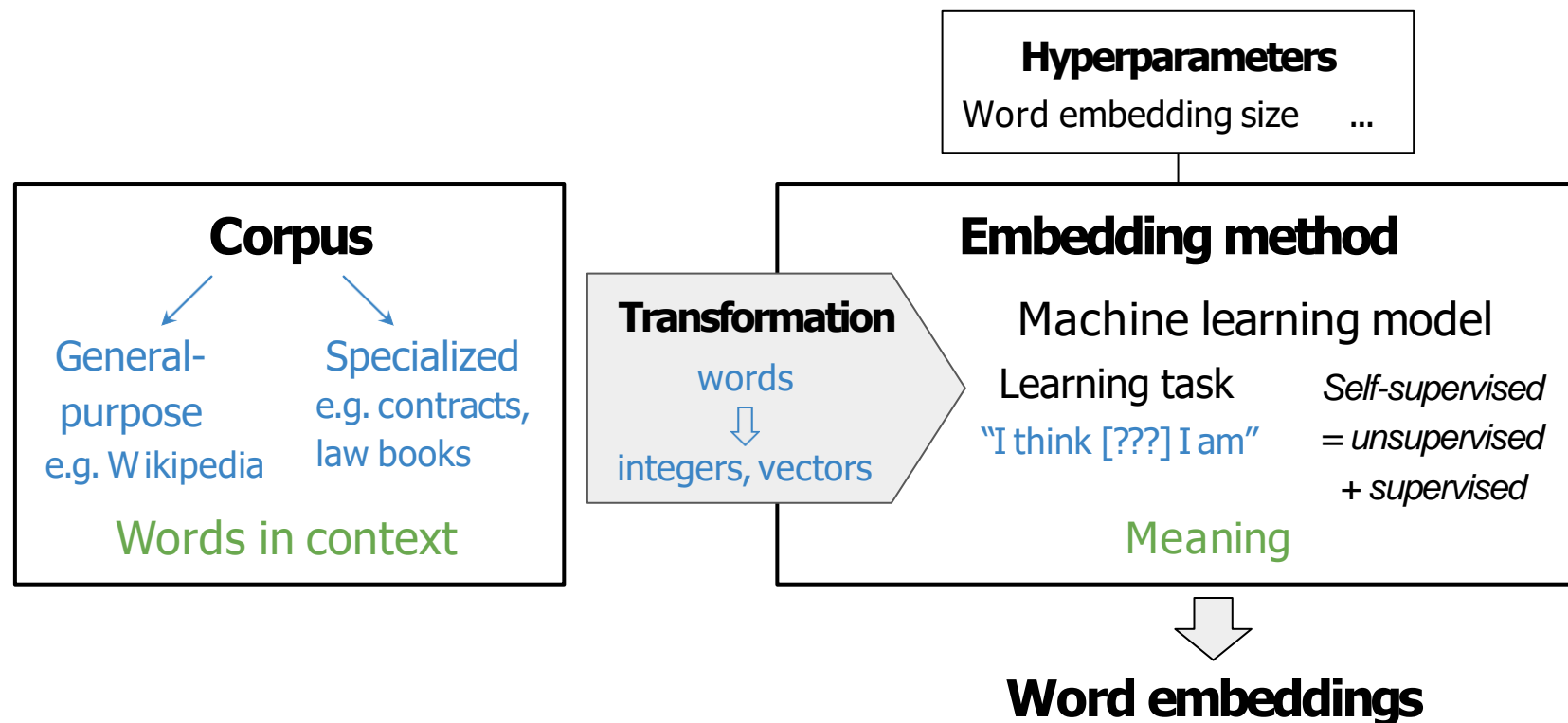| word vectors | |
|---|---|
| one-hot vectors | word embedding vectors |
| | "word vectors" |
| | word embeddings |

integers

# How to Create  Word  Embeddings
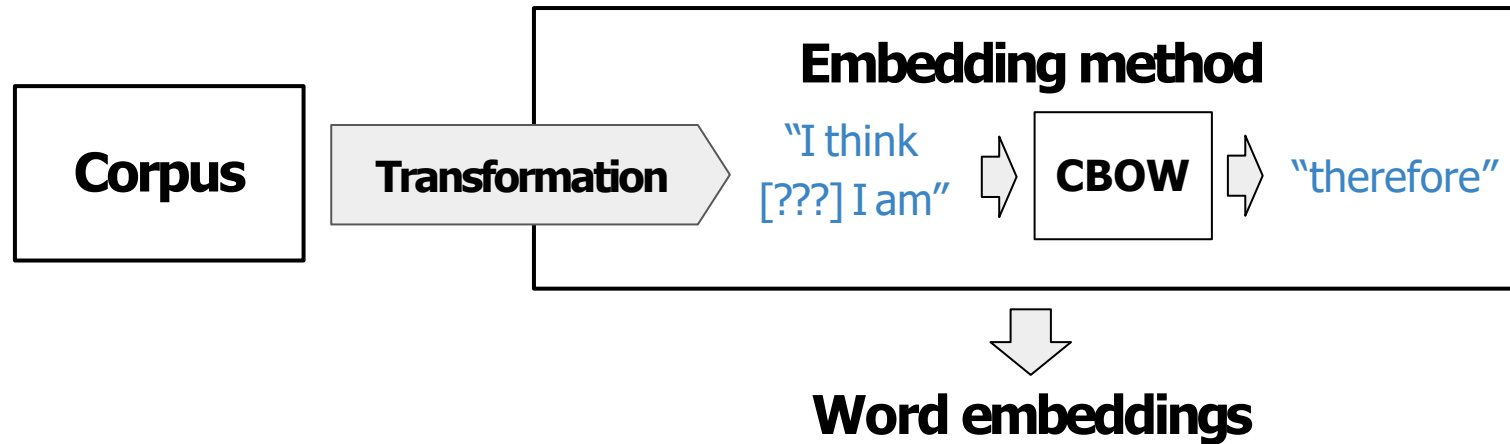
- Word embedding process

# Word embedding methods

- Basic

  - 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

- Deep learning, contextual embeddings

  - BERT (Google, 2018)

  - ELMo (Allen Institute for AI, 2018)
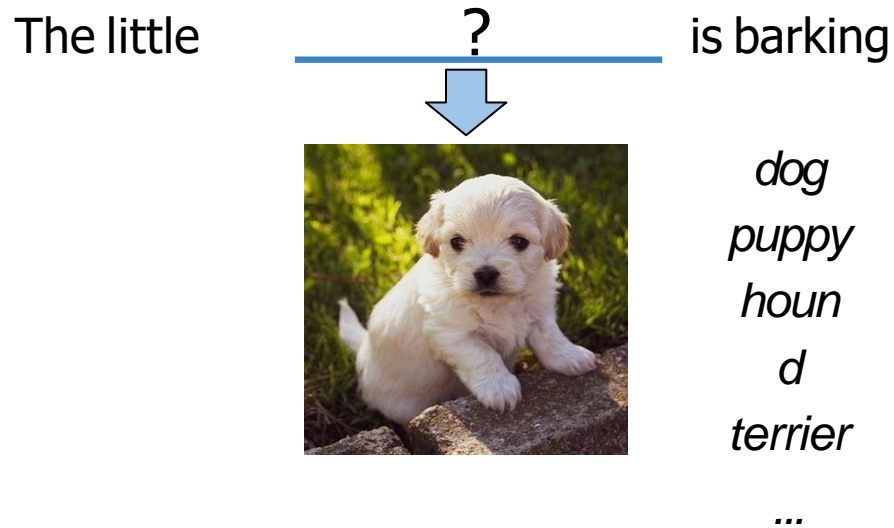
  - GPT-2 (OpenAI, 2018)

Tunable pre-trained models available

# Continuous Bag-of-Words Model

- Continuous bag-of-words word embedding process

# Center word prediction: rationale

The little _____ **?** is barking



_dog_
_puppy_
_houn_
_d_
_terrier_
_..._

Creating a training example

**center word**

I am | happy | because I | am learning

**context words**

**window**

C = 2
context half-size

window size = 5

# From corpus to training

## Corpus

I am happy because I am learning

## Embedding method

### Transformation

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

Context words

→

**CBOW**

→

Predicted center word

↓

## Word embeddings

# Cleaning and  Tokenization

- 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 , libraries

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

emoji

punctuation

number

```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
```

```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
```

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

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

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

# Transforming center words into vectors

Corpus      I am happy because I am learning

Vocabulary    am, because, happy, I, learning

One-hot vector

|          | am | because | happy | I | learning |
|----------|----|---------|-------|---|----------|
| am       | 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

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

$+$ $+$ $+$ $/ 4 =$

I am because I

$$
\begin{pmatrix}
0.25 \\
0.25 \\
0 \\
0.5 \\
0
\end{pmatrix}
$$

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] |

# 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**

$\mathbf{W_1}$ weights

$\mathbf{b_1}$ biases

ReLU

**h**

$\mathbf{W_2}$ weights

$\mathbf{b_2}$ biases

softmax

$\mathbf{\hat{y}}$

Center word vector

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

V

# Dimensions (single input)

Input layer

Hidden layer

Output layer

$W_1$  N x V

$b_1$  N x 1

ReLU

$z_1 = W_1 x + b_1$  N x 1

$x$

V x 1

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

$W_2$  V x N

$b_2$  V x 1

softmax

$z_2 = W_2 h + b_2$  V x 1

$h$

N x 1

$\hat{y} = \text{softmax}(z_2)$  V x 1

$\hat{y}$

V x 1

## Column vectors

$z_1 = W_1 x + b_1$

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

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

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

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

## Row vectors

$z_1 = x W_{1_T} + b_1$

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

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

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

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

# Dimensions (batch input)

$$b_1 \rightarrow B_1 = \begin{pmatrix} b_1 & \dots & b_1 \end{pmatrix} \updownarrow N \quad \textit{broadcasting}$$

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

Input layer      Hidden layer      Output layer

Context words vectors
* matrix

$$X = \begin{pmatrix} x^{(1)} & \dots & x^{(m)} \end{pmatrix} \updownarrow V$$

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

$W_1$    N x V

$B_1$    N x m

ReLU

$W_2$    V x N

$B_2$    V x m

softmax

$Z_1 = W_1 X + B_1$ N x m

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

$Z_2 = W_2 H + B_2$   V x m

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

**X**
V x m

**H**
N x m

**$\hat{Y}$**
V x m

# Dimensions (batch input)

Context words matrix

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

$V$

$m$

Input layer

Hidden layer

Output layer

Predicted center word matrix

$\mathbf{W_1}$

$\mathbf{B_1}$

ReLU

$\mathbf{W_2}$

$\mathbf{B_2}$

softmax

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

$V$

$m$

$\mathbf{X}$

$V \times m$

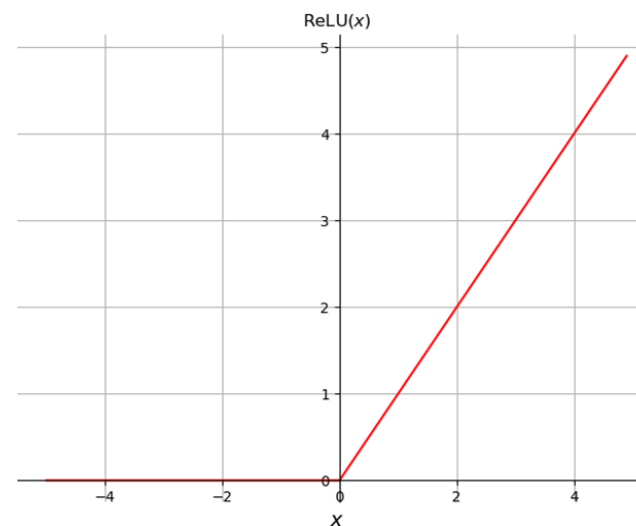$\mathbf{H}$

$N \times m$

$\hat{\mathbf{Y}}$

$V \times m$

# Activation Functions

- Rectified Linear Unit (ReLU)



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

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

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

# Softmax

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

~probabilities

$\Sigma = 1$

Hidden layer     Output layer

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

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

$\mathbf{W_2}$

$\mathbf{b_2}$

softmax

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

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

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

Probabilities of
being center word

# Softmax: example

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

**z**          exp(**z**)          **ŷ** = softmax(**z**)

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

exp →

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

/ Σ →

$$
\begin{pmatrix}
0.083 \\
0.03 \\
0.612 \\
0.225 \\
0.05
\end{pmatrix}
$$

| 0.083 | am |
| 0.03 | because |
| 0.612 | happy |
| 0.225 | I |
| 0.05 | learning |

\*  Predicted center word

Σ=97899

Σ=1

# Cost Function

- Loss

# Cross-entropy loss

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

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

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 |

log →

| $\log(\hat{\mathbf{y}})$ |
|---|
| -2.49 |
| -3.49 |
| -0.49 |
| -1.49 |
| -2.49 |

→

| $\mathbf{y} \odot \log(\hat{\mathbf{y}})$ |
|---|
| 0 |
| 0 |
| -0.49 |
| 0 |
| 0 |

- Σ → $J = 0.49$

# 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

- 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 |

\* J = 4.61



incorrect predictions: penalty

correct predictions: reward

# Training process

- Forward propagation

$$Z_1 = W_1X + B_1 \qquad Z_2 = W_2H + B_2$$
$$H = ReLU(Z_1) \qquad \hat{Y} = softmax(Z_2)$$

# 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

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

Actual center word matrix

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

# Backpropagation and Gradient Descent

- 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

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

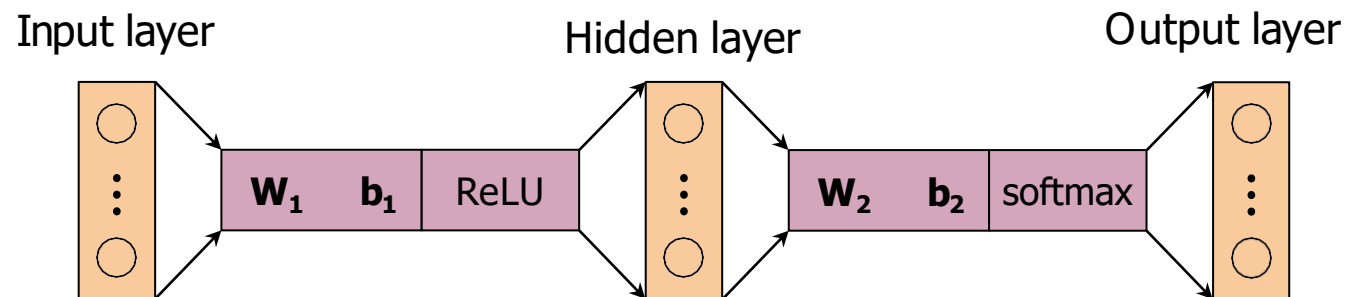- Gradient descent: update weights and biases

# Backpropagation

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

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

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

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

## Gradient descent

Hyperparameter: learning rate $a$

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

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

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

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

# Extracting word embedding vectors: option 1,2,3

# Evaluating Word Embeddings

- Intrinsic evaluation

  - Intrinsic evaluation methods assess how well the word embeddings inherently capture the semantic(meaning) or syntactic(grammar) relationships between the words.

  - Test on semantic analogies

  o Using a clustering algorithm to group similar word embedding vectors, and determining of the cluster's capture related words

  Test relationships between words
  Analogies

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

  Syntactic analogies
  "seen" is to "saw" as "been" is to <?>

  ⚡ Ambiguity
  "wolf" is to "pack" as "bee" is to <?> →   swarm? colony?
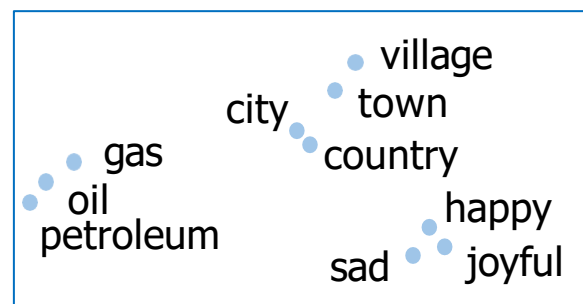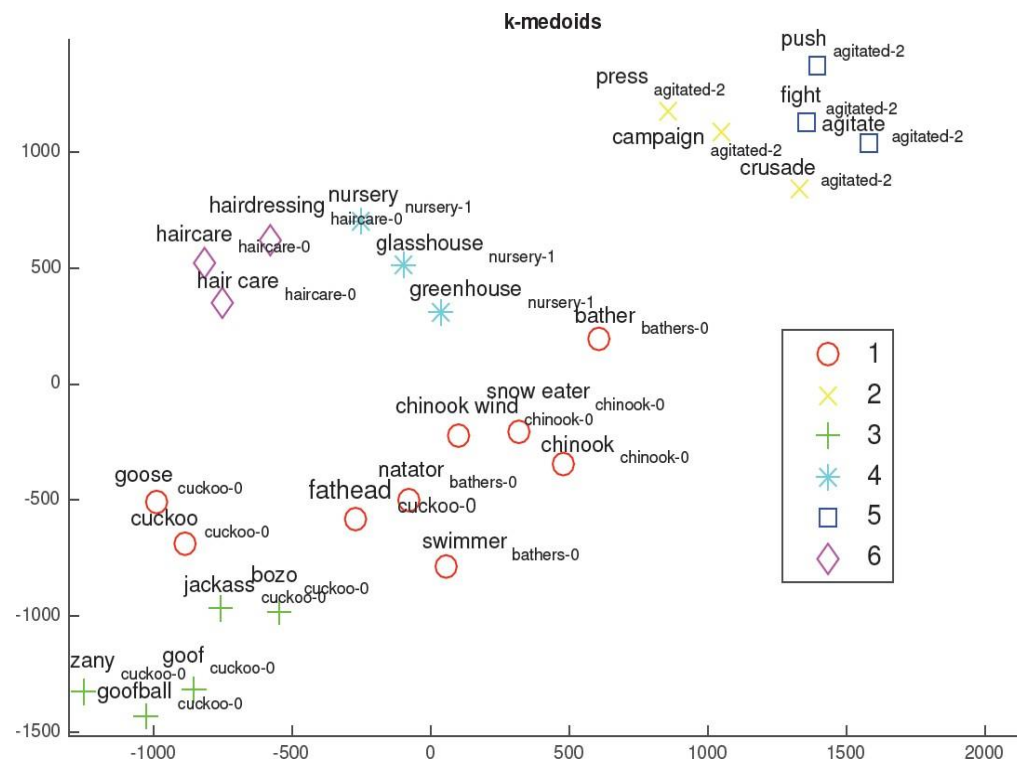
# Intrinsic evaluation

- Test relationships between words

- Analogies

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

# Intrinsic evaluation

- Test relationships between words

  - Analogies

  - Clustering

  - Visualization

# Extrinsic evaluation

o Test the word embeddings to perform an external task, Named Entity recognition, POS tagging

o Evaluate this classifier on the test set with some selected evaluation metric, such as accuracy or the F1 score.

o The evaluation will be more time-consuming than an intrinsic evaluation and more difficult to troubleshoot.

o Test word embeddings on external task

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

o Evaluates actual usefulness of embeddings

    o Time-consuming

    o More difficult to troubleshoot

Named entity

Andrew works at deeplearning.ai

*person*         *organization*