# Machine Learning

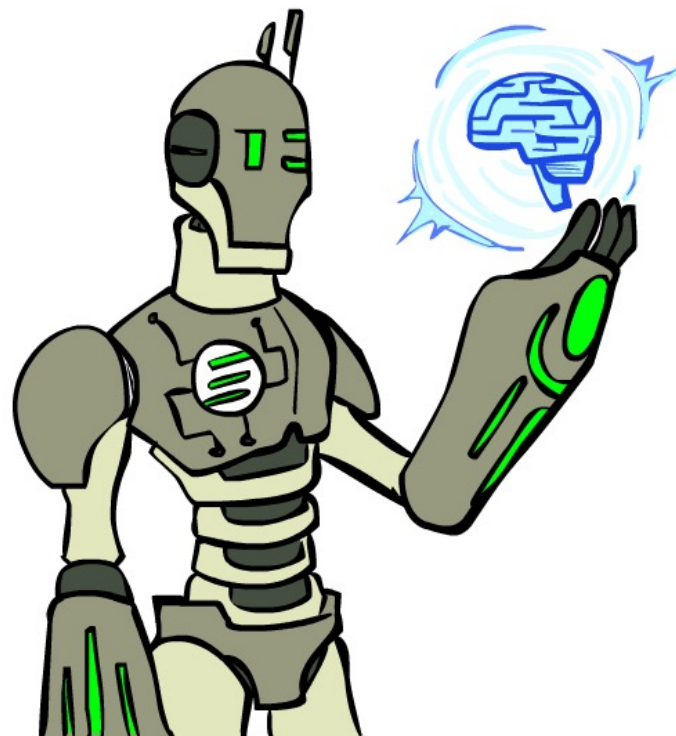## Lecture 9: Application - Word Embedding
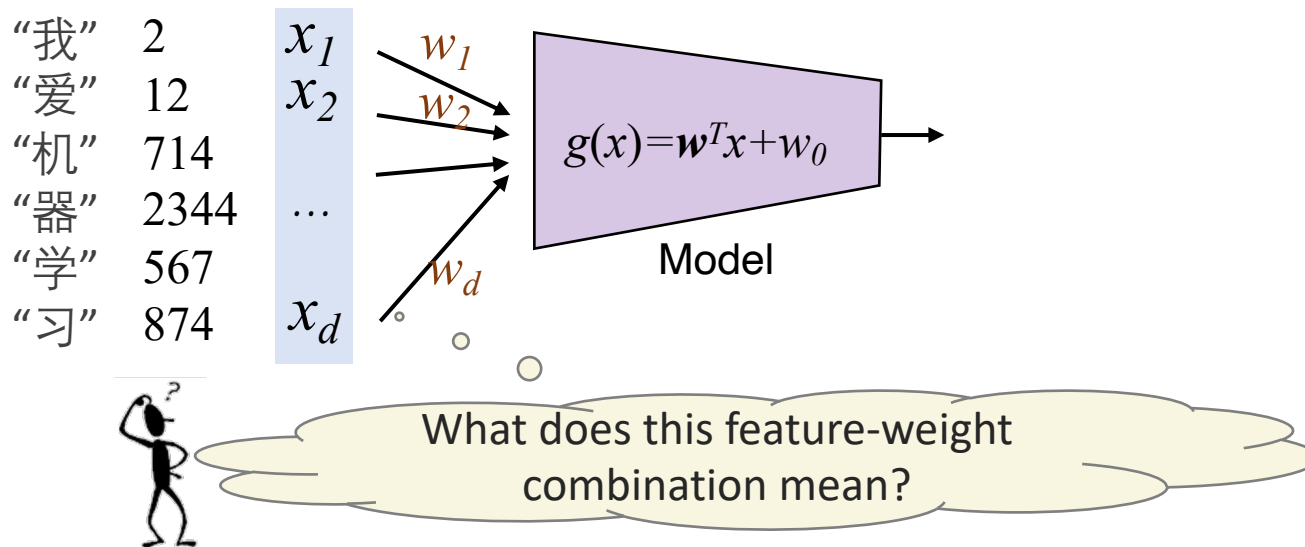
**Fall 2022**

Instructor:  Xiaodong Gu

# **Today**

- Word Embedding (word2vec)

  an application of **neural networks** for distributed representations of words

- Language Model

  a key concept in NLP that is important in the follow-up lectures.

# Before learning word embedding...

- How to represent words in computers?
  - discrete-symbols (ASCII, id, etc.)
  - e.g., "hotel", "conference", "motel", 1203, ...
  - a localist representation

- How to represent words in machine learning models?



"我"  2
"爱"  12
"机"  714
"器"  2344
"学"  567
"习"  874

$x_1$
$x_2$
...
$x_d$

$w_1$
$w_2$
$w_d$

$g(x)=\boldsymbol{w}^T x + w_0$

Model

What does this feature-weight combination mean?

# Words as Vectors

- <mark>Words</mark> can be represented by **one-hot** vectors:

## One-Hot Encoding

> Means one 1, the rest 0s

$$hotel = [\ 1\quad 0\quad 0\quad 0\quad 0\quad 0\quad 0\quad 0\quad 0]$$
$$flower = [\ 0\quad 1\quad 0\quad 0\quad 0\quad 0\quad 0\quad 0\quad 0]$$
$$tree = [\ 0\quad 0\quad 0\quad 1\quad 0\quad 0\quad 0\quad 0\quad 0]$$
$$motel = [\ 0\quad 0\quad 0\quad 0\quad 0\quad 0\quad 0\quad 1\quad 0]$$
$$elephant = [\ 0\quad 0\quad 0\quad 0\quad 0\quad 0\quad 0\quad 0\quad 1]$$

Vector dimensionality = number of words in vocabulary (e.g., **500,000**)

```
           one-hot  vector      word
1.            0
2.            1   0
              (                )
```

How can we <mark>capture relationships (similarity)</mark> between words?

# Problems of One-hot Encoding

Example: in web search, if user searches for "Minhang motel", we would like to match documents containing "Minhang hotel".

**But**

$$\text{motel} = [0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 1\ 0\ 0\ 0\ 0]$$
$$\text{hotel} = [0\ 0\ 0\ 0\ 0\ 0\ 0\ 1\ 0\ 0\ 0\ 0\ 0\ 0\ 0]$$

orthogonal.

There is no natural notion of **similarity** for one-hot vectors!

Solution:

- Use **WordNet**'s list of synonyms?
- But it is well-known to fail badly: incompleteness, etc.

Instead: learn to encode similarity in the vectors themselves

# Distributed Word Representation

- Represent words as <mark>low-dimensional</mark> **dense vectors** that can reflect their semantic similarities. dense

## One-Hot Encoding
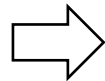
hotel = [ 1  0  0  0  0]

flower = [ 0  1  0  0  0]

business = [ 0  0  1  0  0]

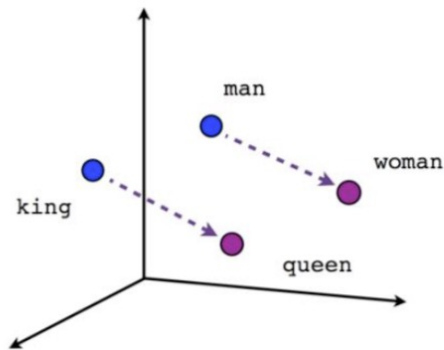motel = [ 0  0  0  1  0]

elephant = [ 0  0  0  0  1]

## Word Embedding



**Note**: word vectors are sometimes called word embeddings or word representations. They are distributed representations.

# Why Word Embeddings?
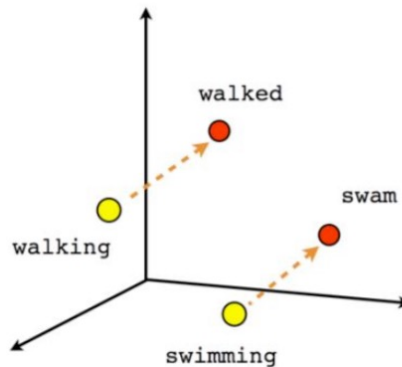
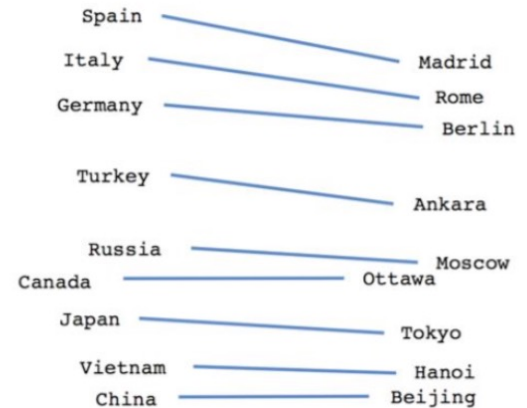- Can capture <mark>the <span style="color:red">rich</span> relational</mark> structure of the lexicon.



Male-Female          Verb tense          Country-Capital

# How to obtain word embeddings?

# How to obtain word embeddings?

- A word can be understood by <mark>its context</mark>

> "A bottle of **tesgüino** is on the table"
> "Everybody likes **tesgüino**"
> "**Tesgüino** makes you drunk"
> "We make **tesgüino** out of corn"



*"You shall know a word by the company it keeps"* (J. R. Firth 1957)

- From context words we can guess **tesgüino** means an alcoholic beverage such as beer.

- Intuition for an algorithm:

<mark>Two words are similar if they have similar word contexts</mark>
idea1

# How to Exploit the Context?

## Counting-based Approach

two words have similar vectors if they frequently co-occur in the same context (sentence, document, etc.).

e.g. Glove:

## Prediction-based Approach

train **neural networks** to predict a word (vector) given its context words (vector).

e.g., word2vec

# Counting based: the ==vector space model==
(            )

# Vector Space Model

- The cornerstone technology in information retrieval.

- **Term-Document Matrix**

  Each cell is the count of word t in document d

  |          | $d_1$ | $d_2$ | $d_3$ | $d_4$ | $d_5$ |
  |----------|-------|-------|-------|-------|-------|
  | ekonomi  | 0     | 1     | 40    | 38    | 1     |
  | pusing   | 4     | 5     | 1     | 3     | 30    |
  | keuangan | 1     | 2     | 30    | 25    | 2     |
  | sakit    | 4     | 6     | 0     | 4     | 25    |
  | Inflasi  | 8     | 1     | 15    | 14    | 1     |

  vector of $d_3$
  = [40, 1, 30, 0, 15]

- Two documents are similar if they have similar vector!

  $d_3$ = [40, 1, 30, 0 15]                    (                    )
  $d_4$ = [38, 3, 25, 4, 14]

# Vector Space Model

- **Term-Document Matrix**

   Each cell is the count of word t in document d

   |            | $d_1$ | $d_2$ | $d_3$ | $d_4$ | $d_5$ |
   |------------|-------|-------|-------|-------|-------|
   | ekonomi    | 0     | 1     | 40    | 38    | 1     |
   | pusing     | 4     | 5     | 1     | 3     | 30    |
   | keuangan   | 1     | 2     | 30    | 25    | 2     |
   | sakit      | 4     | 6     | 0     | 4     | 25    |
   | Inflasi    | 8     | 1     | 15    | 14    | 1     |

   Vector of the word
   "sakit" = [4, 6, 0, 4, 25]

- Two words are similar if they have similar vector!

   pusing = [4, 5, 1, 3, 30]            (            )
   sakit  = [4, 6, 0, 4, 25]

# Vector Space Model

- **Weighting**: in practice, we usually use weights such as TF-IDF, instead of just using <mark>raw counts (only TF)</mark>.

$$\textbf{tf-idf}_{w,\,d} = \text{tf}_{w,\,d} \times \log(N / \text{df}_w)$$

$\text{tf}_{w,\,d}$ = <mark>frequency</mark> of w in d
$\text{df}_w$ = number of documents containing w
$N$ = total number of documents

| | $d_1$ | $d_2$ | $d_3$ | $d_4$ | $d_5$ |
|---|---|---|---|---|---|
| ekonomi | 0 | 1 | 40 | 38 | 1 |
| pusing | 4 | 5 | 1 | 3 | 30 |
| keuangan | 1 | 2 | 30 | 25 | 2 |
| sakit | 4 | 6 | 0 | 4 | 25 |
| Inflasi | 8 | 1 | 15 | 14 | 1 |

$\text{TF(sakit)} = [4, 6, 0, 4, 25]$

$\text{DF(sakit)} = 4$

$N = 5$

$\}\ \text{IDF(sakit)} = \log(5/4)$

$\text{TF-IDF(sakit)} = [\ldots\ldots]$

# Limitations of Vector Space Model

- TF-IDF vectors are
  - long (length |V|= 20,000 to 50,000)
  - sparse (most elements are zero)

  - difficult to use as features in machine learning (more weights to tune)
  - storing explicit counts can be difficult for generalization

# Prediction based: word2vec

# Word2Vec: Overview

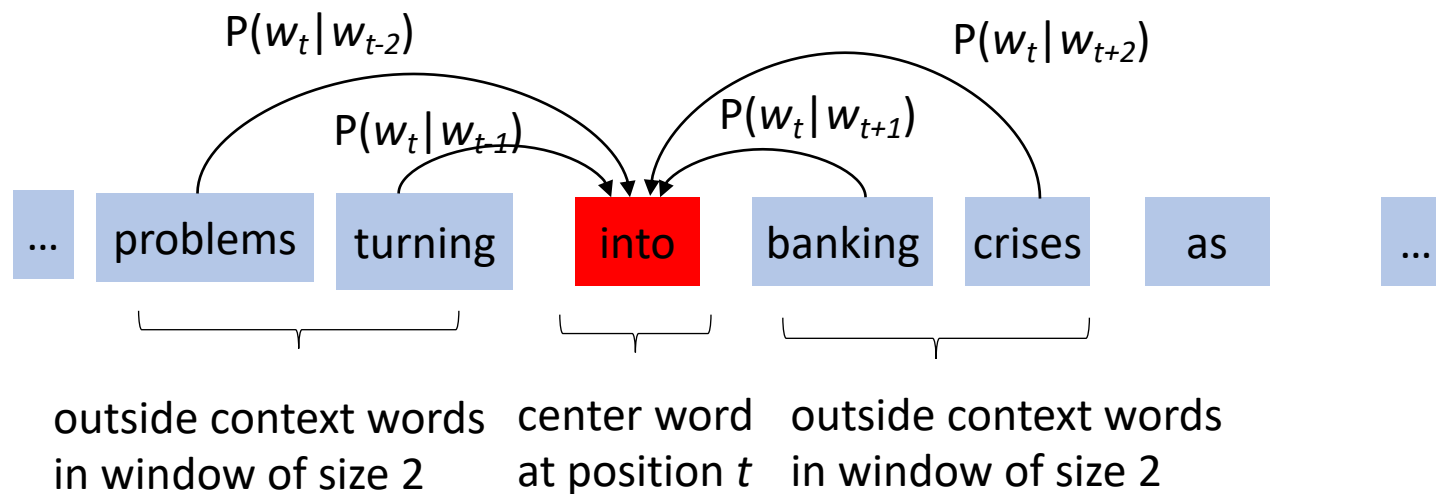- **Word2vec** (Mikolov et al. 2013) is a framework for learning word vectors.

**Idea:**

- we have a large corpus of text.
- every word in a fixed vocabulary is represented by a **learnable vector**;
- go through each position $t$ in the text, which has a center word $c$ and context ("outside") words $o$ ;
- use the similarity of the word vectors for $c$ and $o$ to calculate the probability of $o$ given $c$ (or vice versa) ;
- keep adjusting the word vectors to maximize this probability.

Mikolov T, Sutskever I, Chen K, et al. Distributed representations of words and phrases and their compositionality. NIPS 2013.

# Word2Vec: Overview

Example: <mark>window</mark> and process for computing $P(w_t|w_{t+j})$



$P(w_t|w_{t-2})$

$P(w_t|w_{t+2})$

$P(w_t|w_{t+1})$

$P(w_t|w_{t+1})$

… problems turning **into** banking crises as …

outside context words
in window of size 2

center word
at position $t$

outside context words
in window of size 2

# Word2Vec: Overview

Example: window and process for computing $P(w_t | w_{t+j})$



$P(w_t | w_{t-2})$

$P(w_t | w_{t-1})$

$P(w_t | w_{t+1})$

$P(w_t | w_{t+2})$

... | problems | turning | into | banking | crises | as | ...

outside context words
in window of size 2

center word
at position $t$

outside context words
in window of size 2

# Mikolov's CBOW

- **CBOW**: the distributed representations of context (or surrounding words) are combined to **predict the word in the middle**.

$$P(w_t \mid w_{t-k}, ...w_{t-1}\ w_{t+1}\ ...w_{t+k}) = \text{softmax}\ (\textbf{NN}\ (V(w_{t-k}) + ...+V(w_{t+k})))$$

This can be represented by a **neural network**:

- An input layer which converts each word (one-hot) into a dense vector.
- A projection layer which combines the vectors of input words.
- An output layer which predicts the target word $w_t$ given the combined context vector.



CBOW = Continuous-Bag-of-Words    continuous

the order of words in the context does not influence the projection.

# Model Architecture

- Input Layer: represent any word into a vector. $z_i = Wx_i = W_i$



The weight matrix $W \in R^{|V| \times d}$ is a **lookup table** with each row $W_i$ being the vector for word $w_i$.

one-hot encoding

word vector

### Example

$$[0 \quad 0 \quad 0 \quad 1 \quad 0] \times \begin{bmatrix} 17 & 24 & 1 \\ 23 & 5 & 7 \\ 4 & 6 & 13 \\ 10 & 12 & 19 \\ 11 & 18 & 25 \end{bmatrix} = [10 \quad 12 \quad 19]$$
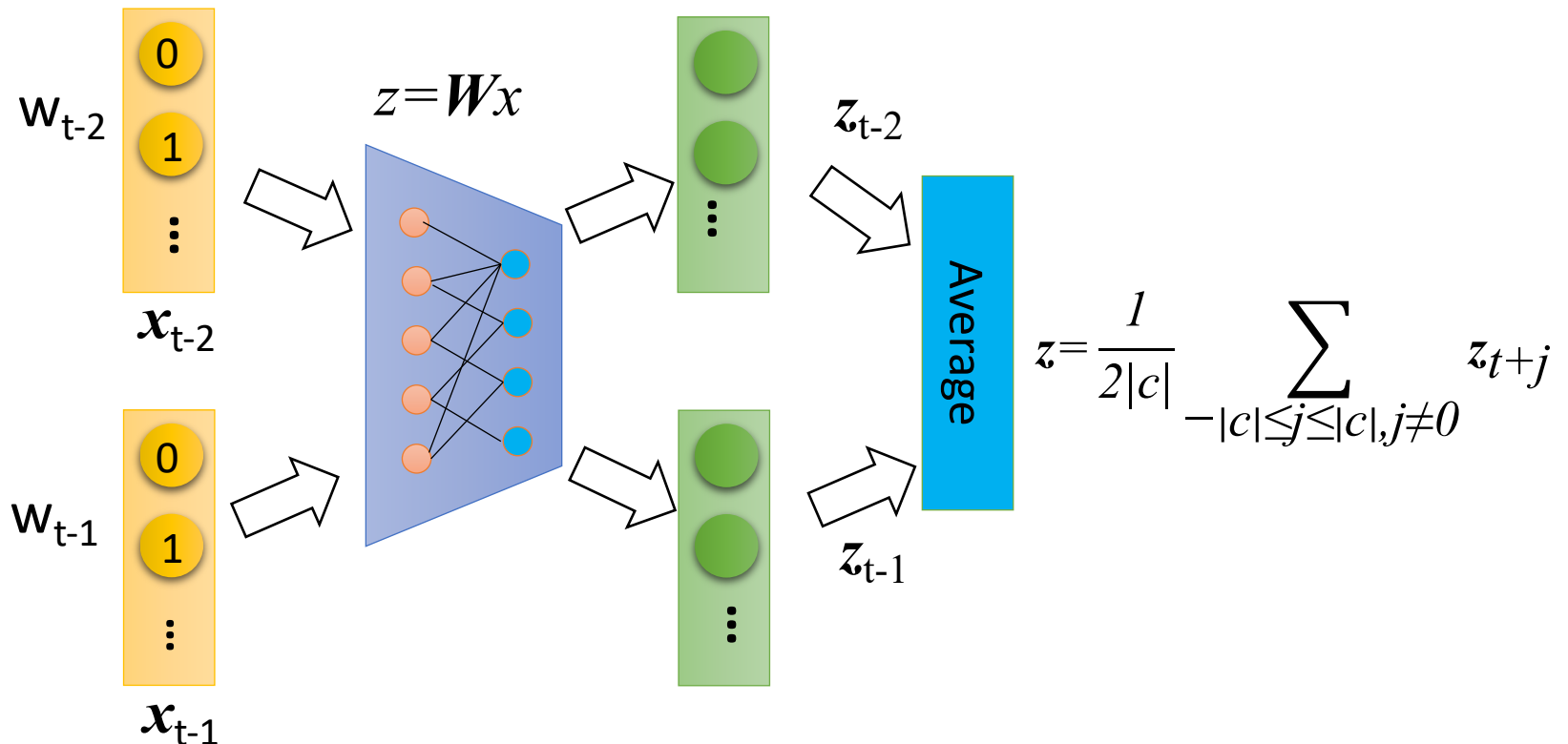
"motel"

motel

hotel

# Architecture

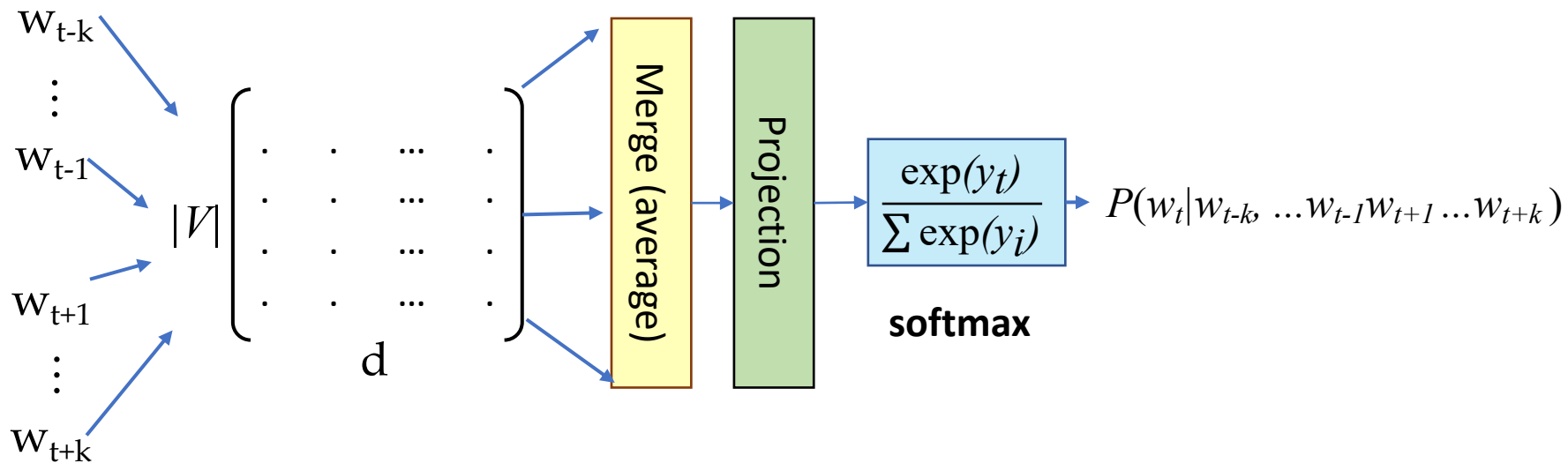- Projection Layer: combining context vectors <mark>into one vector</mark>.



$z=\boldsymbol{W}x$

$\boldsymbol{z}_{t-2}$

$\boldsymbol{z}_{t-1}$

$\mathsf{w}_{t-2}$

$\boldsymbol{x}_{t-2}$

$\mathsf{w}_{t-1}$

$\boldsymbol{x}_{t-1}$

Average

$$z=\frac{1}{2|c|}\sum_{-|c|\leq j\leq|c|, j\neq 0} \boldsymbol{z}_{t+j}$$

# Architecture

- Output Layer: predicts the probability of the target word.

$$P\left(w_t \middle| w_{t-|c|}, \ldots, w_{t-1}, w_{t+1}, \ldots, w_{t+|c|}\right) = \frac{\exp(y_t)}{\sum_{i=1}^{|V|} \exp(y_i)}$$



**Softmax**

$y = \boldsymbol{U}z$

$z_1$
$z_2$

$p(w_t | w_{t+j})$

t

probability of each
word being the target

# Architecture: the big picture

$$w_{t-k}$$

$\vdots$

$$w_{t-1}$$

$$w_{t+1}$$

$\vdots$

$$w_{t+k}$$

$|V| \begin{pmatrix} \cdot & \cdot & \cdots & \cdot \\ \cdot & \cdot & \cdots & \cdot \\ \cdot & \cdot & \cdots & \cdot \\ \cdot & \cdot & \cdots & \cdot \end{pmatrix}$

d

Merge (average)

Projection

$$\frac{\exp(y_t)}{\sum \exp(y_i)}$$

**softmax**

$$P(w_t | w_{t-k}, \ldots w_{t-1} w_{t+1} \ldots w_{t+k})$$

# Training

- Given $D = \{w_1, w_2, \ldots, w_N\}$, minimize the ==negative log likelihood== (NLL) loss function:

$$L(W, U \mid D) = -\frac{1}{N} \sum_{t=1}^{N} \log p(w_t \mid w_{t-k}, \ldots, w_{t-1}, w_{t+1}, \ldots, w_{t+k})$$

using gradient descend.

# Other Models

- **CBOW**    predict the word ==given its context==



- **Skip-gram**    predict the context ==given a word==

# The Skip-Gram Model

- We seek a model for P(<mark>$w_{t+j} | w_t$</mark>).

$$P\left(w_{t+j} | w_t\right) = \frac{\exp(y_{t+j})}{\Sigma_{i=1}^{|V|} \exp(y_i)}$$

$$y = \boldsymbol{U}z$$

$$z = \boldsymbol{W}x$$

INPUT    PROJECTION    OUTPUT

w(t)

w(t-2)
w(t-1)
w(t+1)
w(t+2)

$$L(W, U | \chi) = -\frac{1}{N} \sum_{t=1}^{N} \sum_{-c \leq j \leq c, j \neq 0} \log p(w_{t+j} | w_t)$$

# Example

- Word2vec maximizes an objective function by ==putting similar words nearby in space==.



Fu, Ruiji, et al. "Learning semantic hierarchies via word embeddings." *ACL* 2014.

# The Word Analogy Task <sup>analogy</sup>

- Word Analogy:

  | a:b :: c:? |
  |---|

  man:woman :: king:?

  **Examples**
  - Man is to Woman as King is to___ ?
  - Good is to Best as Smart is to___?
  - China is to Beijing as America is to__?

  How to find d ?

  $$d = \arg\max_i \frac{(x_b - x_a + x_c)^T x_i}{||x_b - x_a + x_c||}$$

- It turns out that word2vec is good for such an analogy task.

  $V_{king} - V_{man} + V_{woman} = V_{queen}$



king

woman

man

Mikolov, T., et al. Distributed Representations of Words and Phrases and Their Compositionality. NIPS 2013.
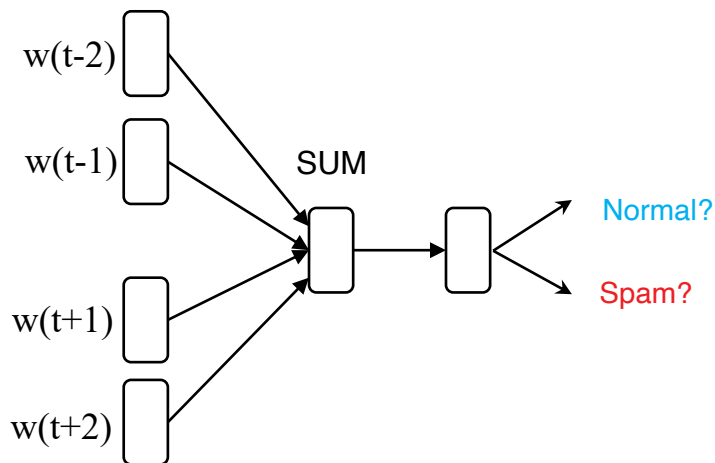
# Example

# Thinking

softmax

- What kind of machine learning is word2vec?

- What limitations does word2vec have if we use it for text classification ?

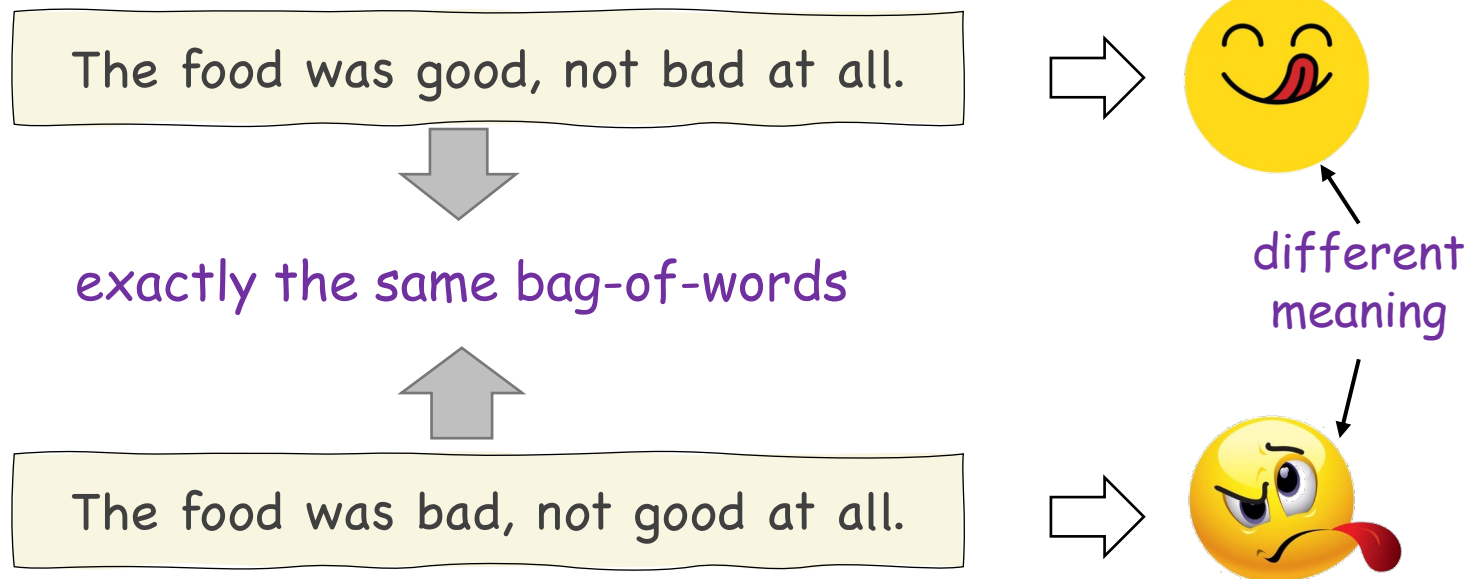# Beyond Bag-of-Words

**Order Matters!**

- To understand the meaning of a sentence, the **order** of the words can not be ignored.

| The food was good, not bad at all. |

exactly the same bag-of-words

| The food was bad, not good at all. |

different meaning

*Paragraph Vector*: Le, Quoc, and Tomas Mikolov. "Distributed Representations of Sentences and Documents." ICML, 2014
*Seq2seq Auto-encoder*: Li, Jiwei, Minh-Thang Luong, and Dan Jurafsky. "A hierarchical neural autoencoder for paragraphs and documents." arXiv preprint, 2015
*Skip Thought*: Ryan Kiros, Yukun Zhu, Ruslan Salakhutdinov, Richard S. Zemel, Antonio Torralba, Raquel Urtasun, Sanja Fidler, "Skip-Thought Vectors" arXiv preprint, 2015.

# Beyond Bag-of-Words

- **Language Models**
  Recurrent Neural Networks, Transformers,…

- **Pretrained Language Models**
  BERT, GPT-2, …

To appear in future lectures..

# Language Models

- A probabilistic model of how likely a given string appear in a given "language".

- For any sequence $x = (w_1, w_2, \ldots, w_N)$, a **language model** can be defined as:

string

language
string

$$p(x) = p(w_1, w_2, \ldots, w_N)$$

**Example:**

$P_1 = P(\text{"我爱机器学习"})$
$P_2 = P(\text{"我爱学习机器"})$
$P_3 = P(\text{"机器我爱学习"})$
$P_4 = P(\text{"爱我机学习器"})$

Chinese: $P_1 > P_2 > P_3 > P_4$

- Applications:

  message suggestion; document generation; spelling correction; machine translation; speech recognition;…

我爱机器学 __?__

# Language Model

- What is the probability of $P(w_1, \ldots, w_N)$?

> p(我爱机器学习) = ?

- Chain Rule:

$$p(w_1, \ldots, w_N) = p(w_1)p(w_2|w_1)\ldots,p(w_N|w_1,\ldots,w_{N-1})$$

> p(我爱机器学习) = p(我)p(爱|我)p(机|我爱)p(器|我爱机)p(学|我爱机器)p(习|我爱机器学)

- Markov Assumption: (only consider the last n-1 words)

$$p(w_i|w_1,\ldots,w_{i-1}) = p(w_i|w_{i-n+1},\ldots,w_{i-1})$$

> p(习|我爱机器学) ≈ p(习|机器学) ≈ p(习|学)

---

# Bigram Language Model

So that's what we get for <mark>n=2</mark>:

$$p(w) = p(w_1)p(w_2|w_1)\ldots,p(w_N|w_{N-1})$$

$$1/18 \times 1/8 \times 1/120 \times 1/4 \times 1/420 \times 1/2$$

p(我爱机器学习) = p(我)p(爱|我)p(机|爱)p(器|机)p(学|器)p(习|学)

Q: How to estimate these probabilities?

A: Straightforward counting.

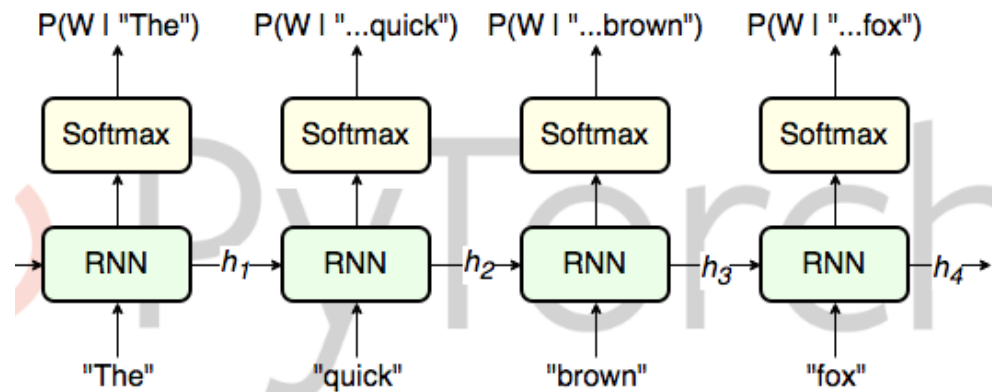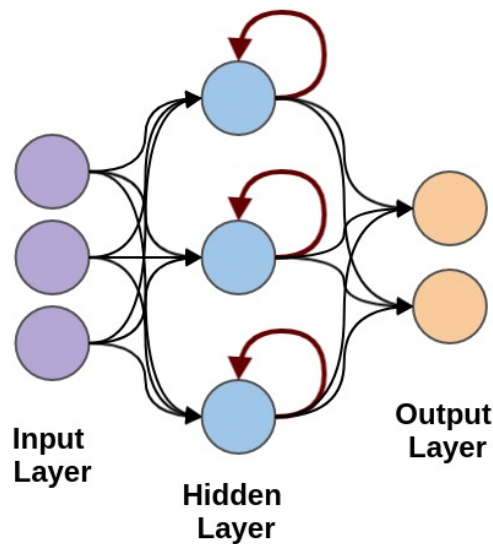Q: But remember in word embedding, counting discrete words have many drawbacks?

A: Don't worry, we have neural networks as language models.

# What's Next?

## Recurrent Neural Networks

- A deep neural network for sequence (language) modeling.