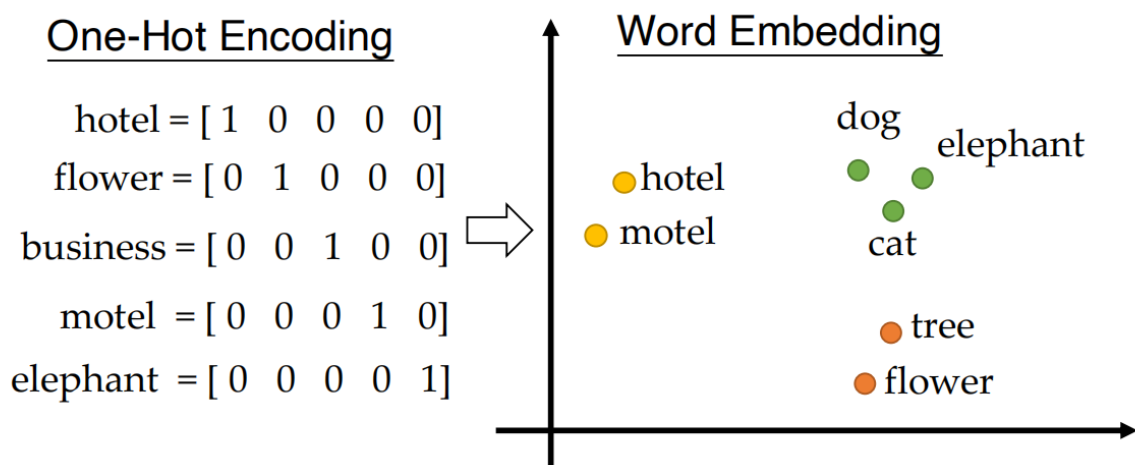# ⛅ ch9: Word Embedding

## Word Embedding

### overview(P6)

- Represent words as low-dimensional **dense vectors** that can reflect their semantic similarities. *由sparse变dense*



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

### Why Word Embeddings?

> Can capture the rich relational structure of the lexicon.(抓住词之间的语义关系)

- **Two words are similar if they have similar word contexts**

## Models

### Counting based: the vector space model(P10-14)

- The cornerstone technology in information retrieval.

信息检索中的重要算法

- **Term-Document Matrix**

  Each cell is the count of word t in document d

  在文档 $d_4$ 中出现的次数

|  | $d_1$ | $d_2$ | $d_3$ | $d_4$ | $d_5$ |
|---|---|---|---|---|---|
| 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]

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

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

惩罚项

tf$_{w, d}$ = frequency of w in d
df$_w$ = number of documents containing w → 如果在所有文档当中都
N = total number of documents    出现, 说明重要性一般

|  | $d_1$ | $d_2$ | $d_3$ | $d_4$ | $d_5$ |
|---|---|---|---|---|---|
| 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 |

TF(sakit) = [4, 6, 0, 4, 25]

DF(sakit) = 4
N = 5 } IDF(sakit) = log(5/4)

TF-IDF(sakit) = [......]

- **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

在机器学习中很难作为特征使用（更多的权重来调整）
存储显式计数很难泛化

## Prediction based: word2vec(P16-26)

- overview(P17)

- **Word2vec** (Mikolov et al. 2013) is a framework for learning word vectors.
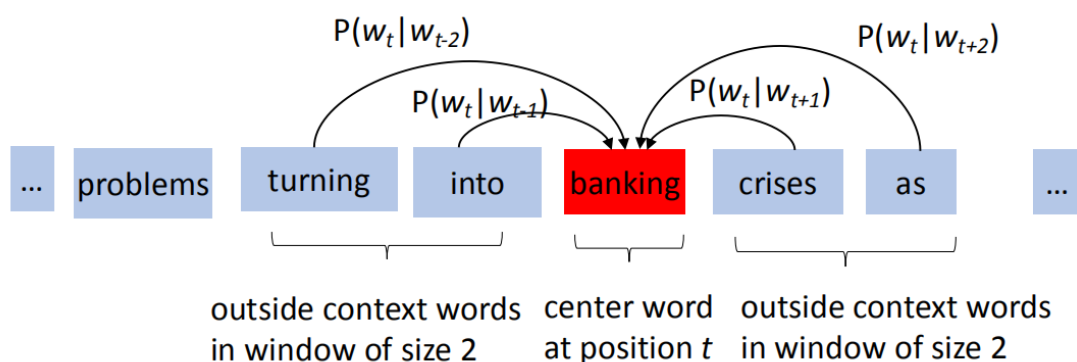
Idea:
- we have a large corpus of text.
- every word in a fixed vocabulary is represented by a vector; *learnable，可变的*
- go through each position $t$ in the text, which has a center word $c$ and context ("outside") words $o$; *保证word和neighbor word共同*
- use the similarity of the word vectors for $c$ and $o$ to calculate the probability of $o$ given $c$ (or vice versa); *出现的概率尽可能高*
- keep adjusting the word vectors to maximize this probability.

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
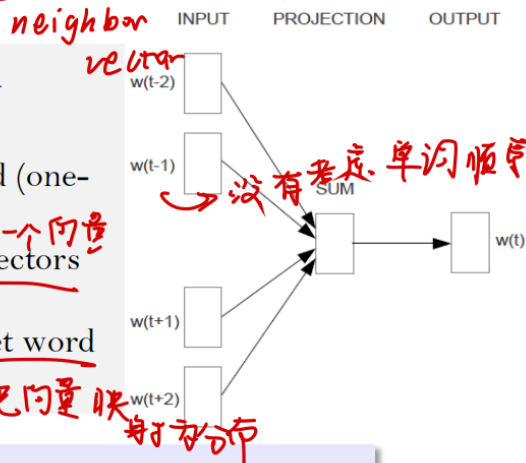
- **Mikolov's CBOW(P20)**

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

$$P(w_t|w_{t-k}, ...w_{t-1} \; w_{t+1} ...w_{t+k}) = NN \; (V(w_{t-k}) + ...+V(w_{t+k}))$$

*neighbor vector*

INPUT     PROJECTION     OUTPUT

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

单词→向量
- An input layer which converts each word (one-hot) into a dense vector. 变成一个向量
- 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. 把向量映射到输出
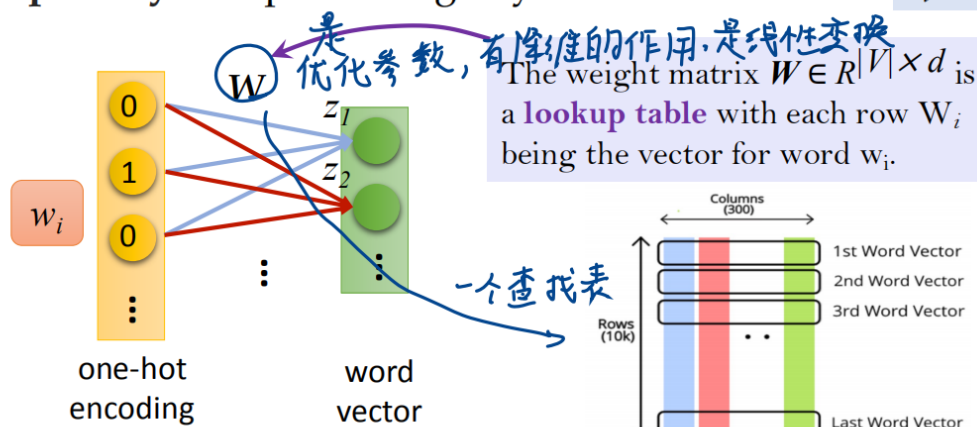
→ 没有考虑 单词顺序

w(t-2)

w(t-1)

SUM

w(t+1)

w(t+2)

w(t)

CBOW = Continuous-Bag-of-Words
the order of words in the context does not influence the projection.

- **Architecture(P21-24)**

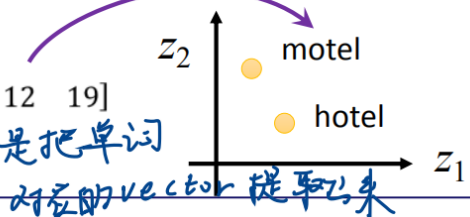- **Input Layer:** representing any word into a vector.   $z_i = Wx_i = W$
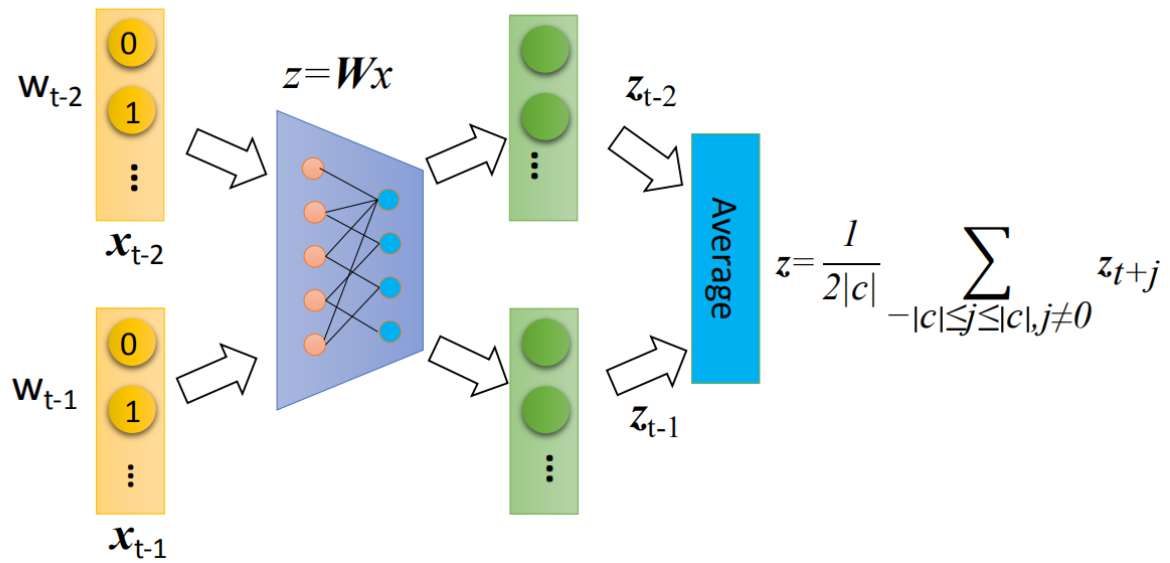
是优化参数，有降维的作用，是线性变换

$W$

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

$z_1$

$z_2$

$w_i$

一个查找表

one-hot encoding     word vector

Columns (300)

Rows (10k)

1st Word Vector
2nd Word Vector
3rd Word Vector

Last Word Vector

Example

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

"motel"

本质上是把单词对应的 vector 提取出来

$z_2$   motel

hotel

$z_1$

- **Projection Layer**: combining context vectors into one vector.



$$z = Wx$$

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

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

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

要预测的层的概率分布

Softmax →转化到词汇表的维度

$$y = Uz$$

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

probability of each word being the target



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

softmax

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

- **Training(P25)**

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

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

与 交叉熵等价

using gradient descend.

## The Skip-Gram Model（P27）
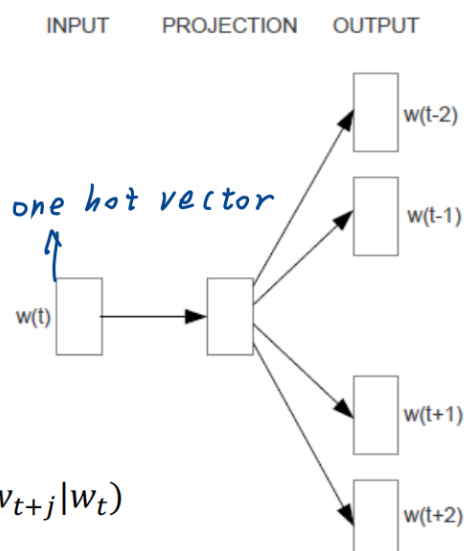
- We seek a model for $P(w_{t+j} \mid w_t)$.

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

$$y = Uz$$

$$z = Wx$$

想 要得到的

别助词
出现的概率

INPUT　　PROJECTION　　OUTPUT

one hot vector



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

## The Word Analogy Task(P29-30)

- Word Analogy: 单词类推

$$a:b :: c:?$$

man:woman :: king:?

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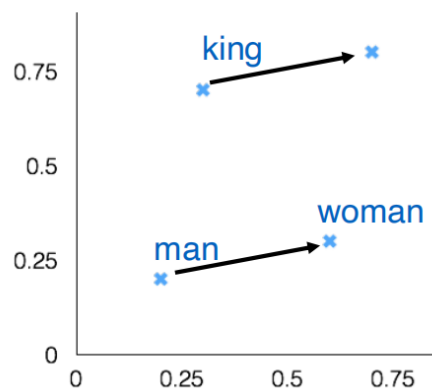- Man is to Woman as King is to___?
- Good is to Best as Smart is to___?
- China is to Beijing as America is to___?

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.

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



## Language Models(P33-35)

- A probabilistic model of how likely a given string appear in a given "language". 评估一个string在语言中出现的概率
- For a given sequence $x = (w_1, w_2, ..., w_N)$. A **language model** can be defined as:

$$p(x) = p(w_1, w_2, ..., w_N)$$

概率越高，对语言的掌握程度越高

**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;…

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

$$p(我爱机器学习) = ?$$

- Chain Rule:

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

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

不是一个词生成时，只看前面 n-1 个

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

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

$$p(习|我爱机器学) \approx p(习|机器学) \approx p(习|学)$$

So that's what we get for n=2: 只看前面一个

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

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

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