

Representation Learning & Transfer Learning

张江

北京师范大学系统科学学院教授

集智俱乐部、集智学园创始人

集智研究中心理事长

What is deep learning?



Deep
Learning

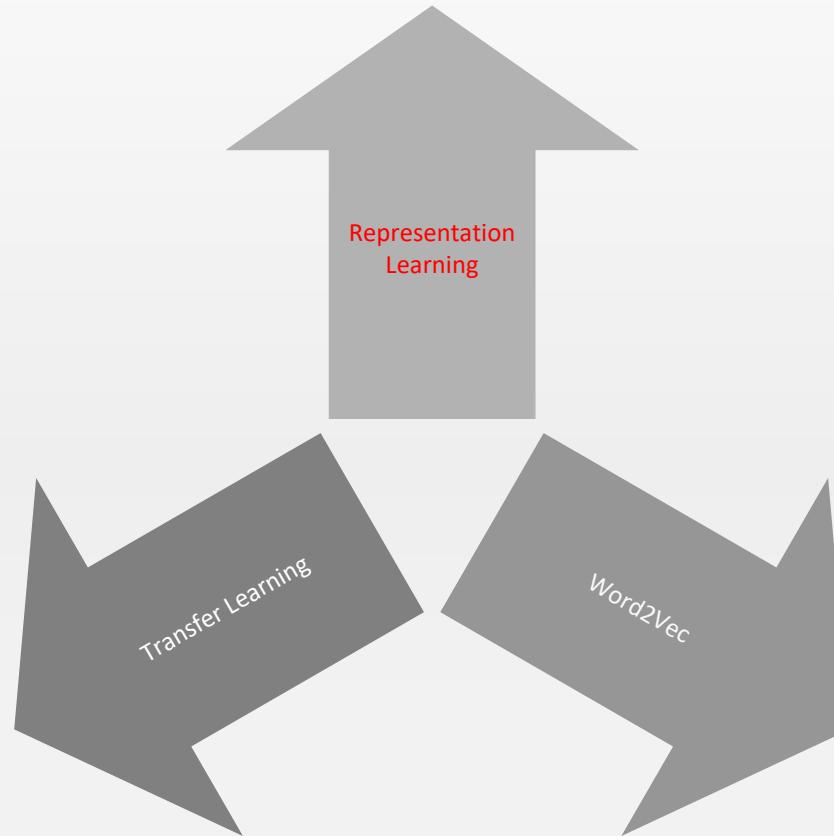


~~Deep?~~

Automatic Differentiation

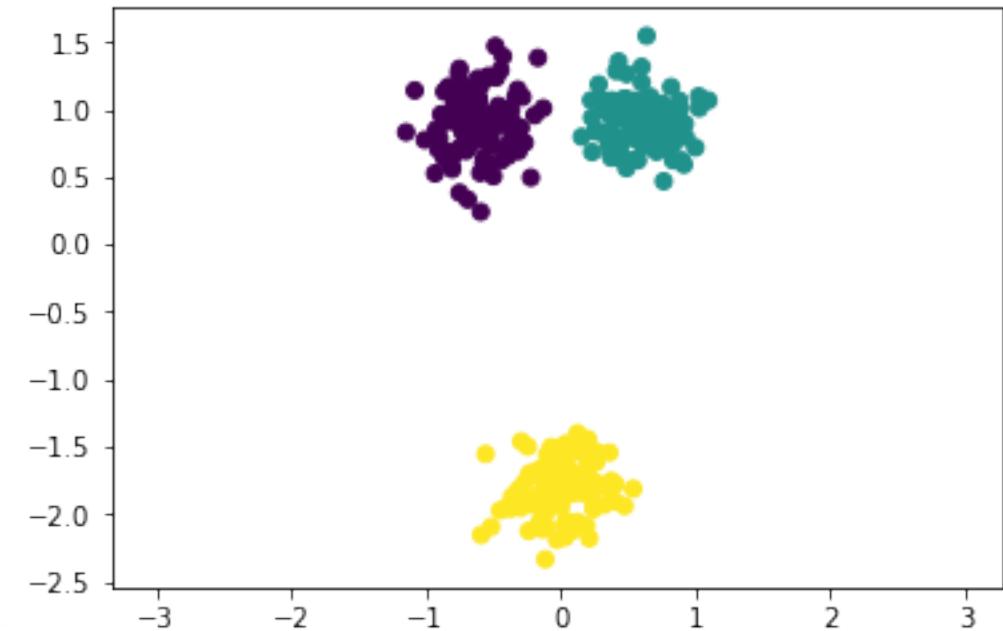
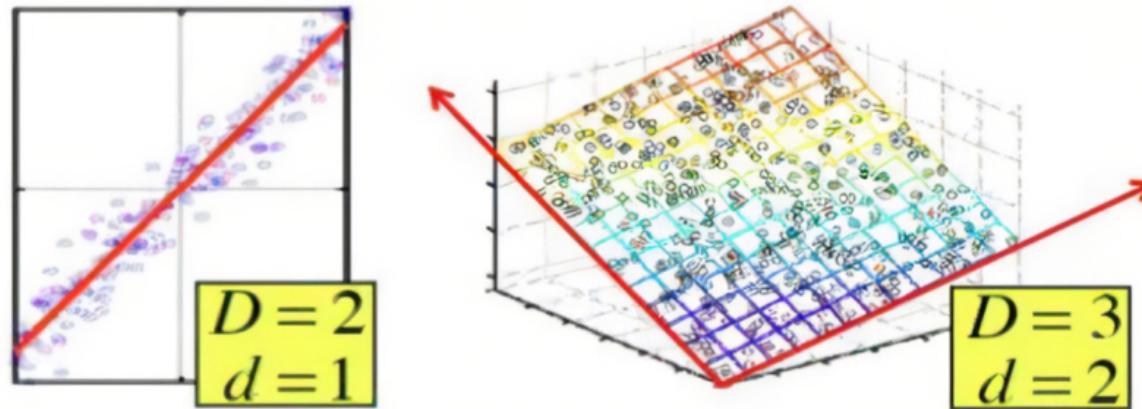
Representation
Learning

Outline



Problem I – Dimension Reduction

Dimensionality Reduction



Clustering

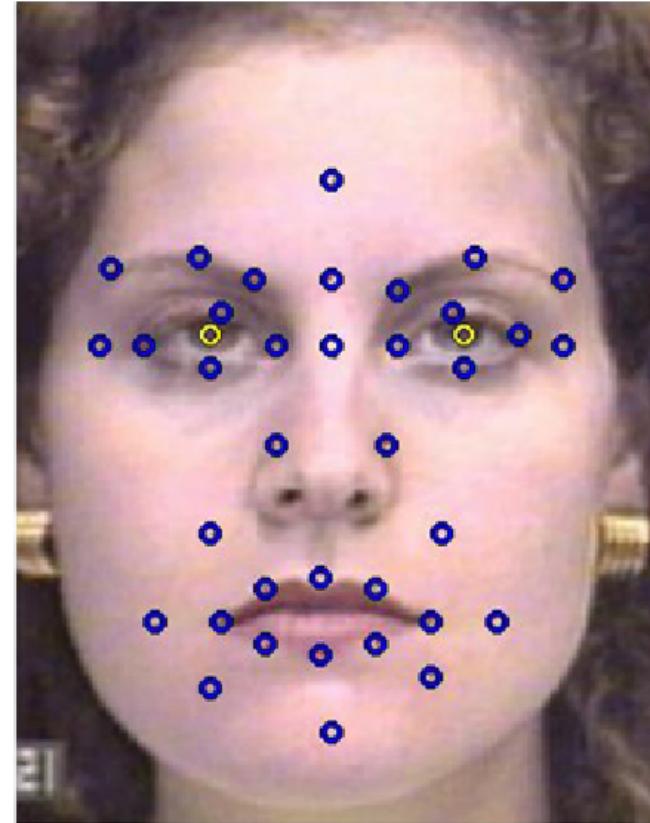
Problem II – Feature Extraction

Multi-State Models for Geometric Feature Extraction



(a)

- One*** State for Brow and Cheek
- Two*** States for Eye
- Two*** States for Furrows
- Three*** States for Lip



(b)

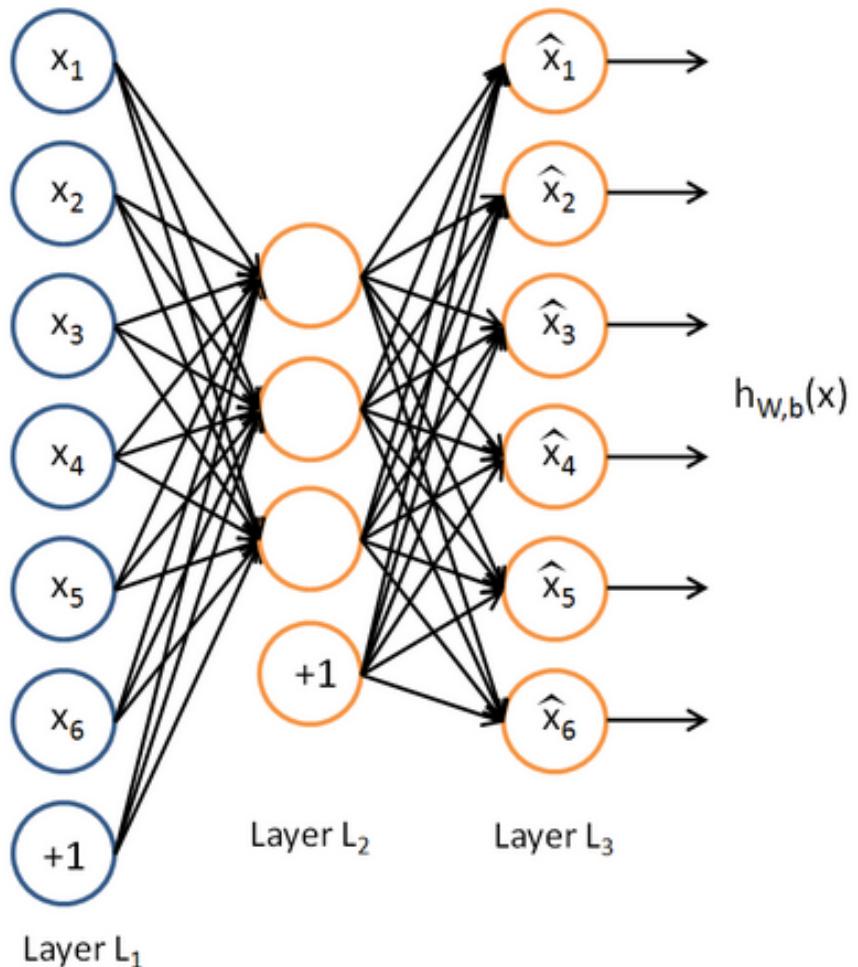
Problem III – Representation of Complex Objects



Categorical variables are a natural choice for representing discrete structure in the world. However, stochastic neural networks rarely use categorical latent variables due to the inability to backpropagate through samples. In this work, we present an efficient gradient estimator that replaces the non-differentiable sample from a categorical distribution with a differentiable sample from a novel Gumbel-Softmax distribution. This distribution has the essential property that it can be smoothly annealed into a categorical distribution. We show that our Gumbel-Softmax estimator outperforms state-of-the-art gradient estimators on structured output prediction and unsupervised generative modeling tasks with categorical latent variables, and enables large speedups on semi-supervised classification.



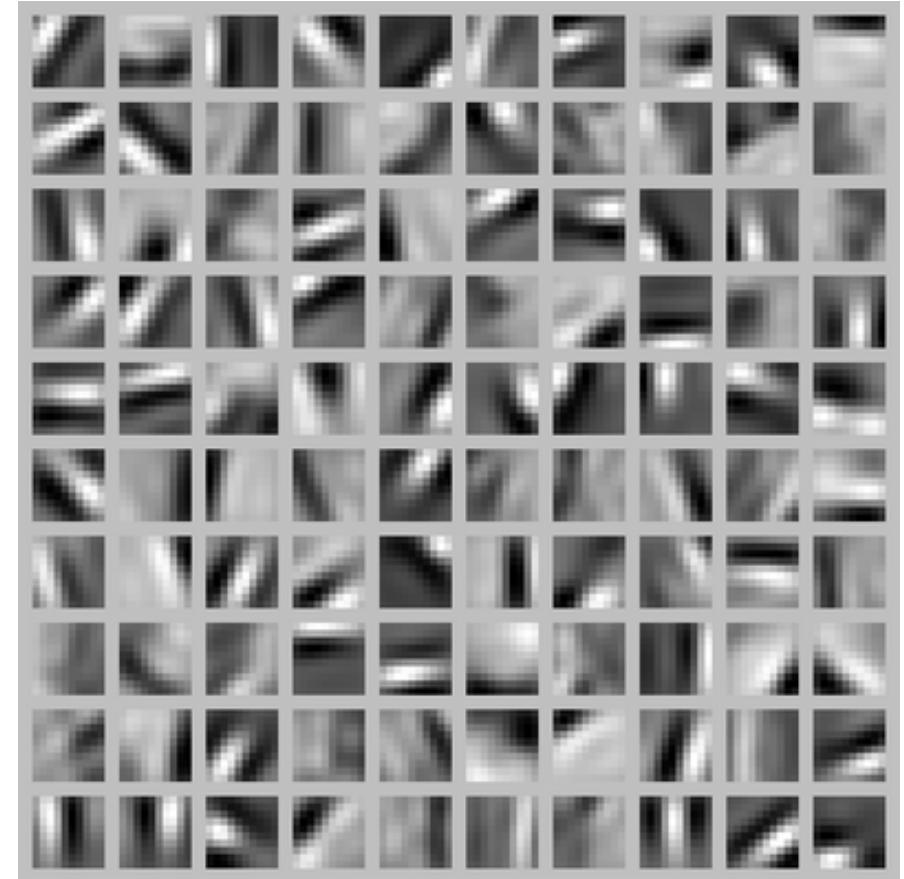
Autoencoder



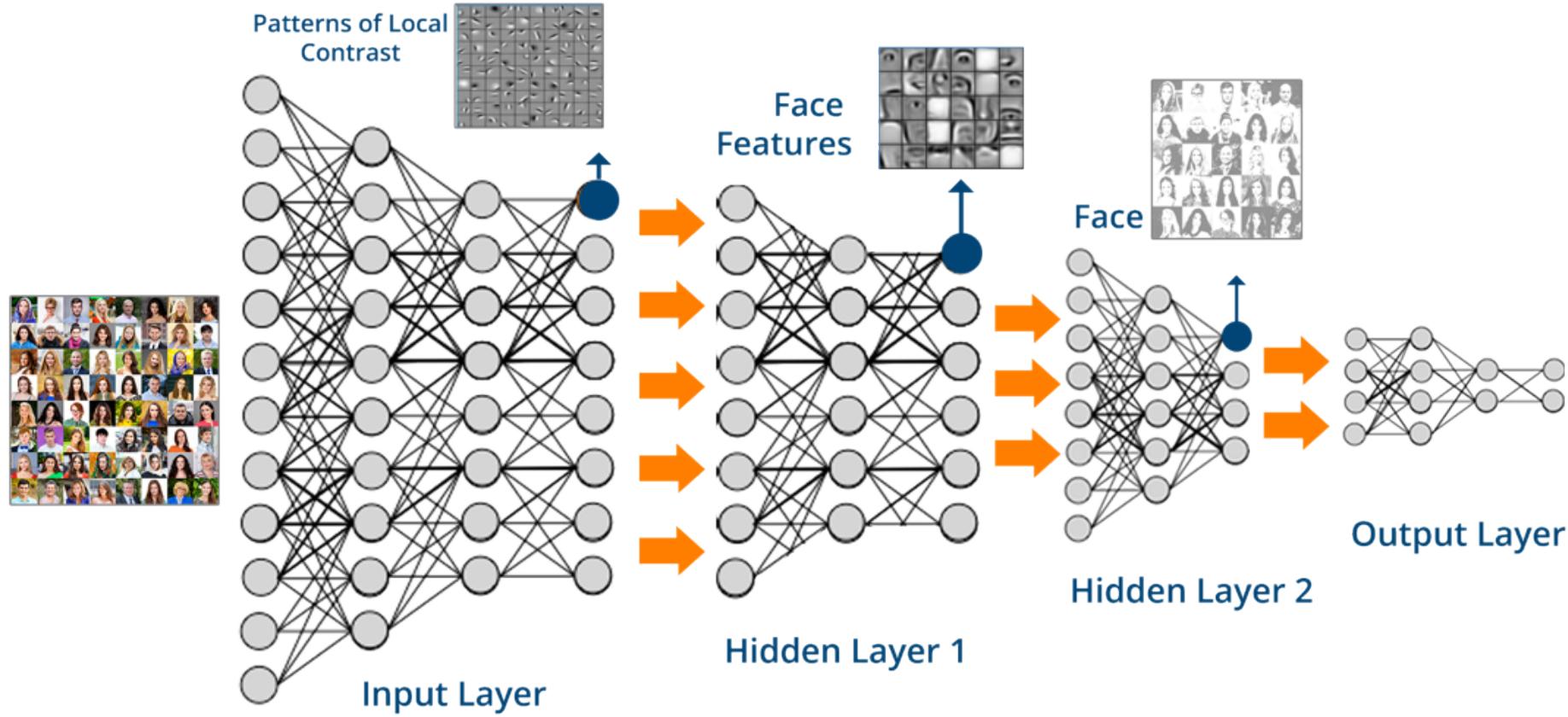
- Unsupervised learning
- Feed-back
 - Reconstruct input $h_\theta(x) \approx x$
- BP Learning
- Constraints

Autoencoder

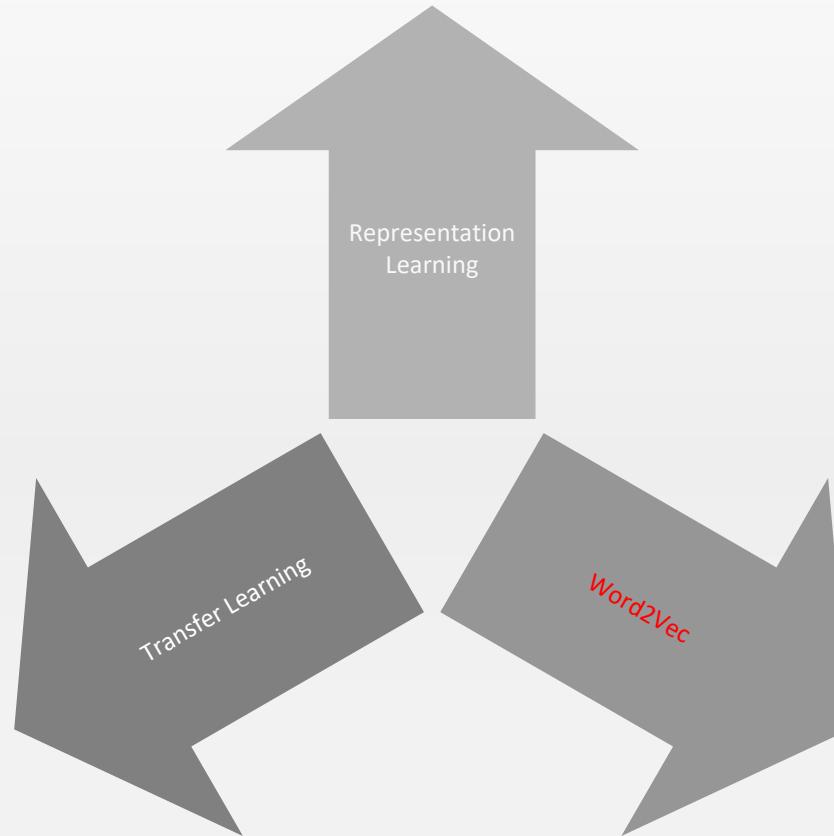
- Constraint activation
- Describe the input data with minimum features.
- When will one unit is activated?
 - when w parallels with x_i



Features of Depth



Outline



NLP is a field of computer science and linguistics concerned with the interaction between computers and human language. It involves the study of natural languages and the development of algorithms and systems that can process, analyze, and generate language. NLP has many applications, including machine translation, speech recognition, text mining, and sentiment analysis.

The field of NLP has been around since the 1950s, but it really began to take off in the 1980s with the development of statistical models. These models use large amounts of data to train computers to recognize patterns in language. This has led to significant advances in areas like POS tagging, named entity recognition, and dependency parsing.

One of the main challenges in NLP is dealing with the complexity of human language. Languages are highly structured and have many nuances that are difficult to capture with simple rules. This is why many NLP systems use machine learning techniques to learn from examples rather than being programmed by hand. However, this also means that NLP systems can be less reliable than rule-based systems, especially when it comes to tasks like sentiment analysis or text generation.

Another challenge is dealing with the diversity of languages in the world. There are thousands of languages spoken around the globe, each with its own unique grammar and vocabulary. This makes it difficult to develop a single system that can handle all of them. Instead, most NLP systems focus on specific languages or language families, such as English or German.

Despite these challenges, NLP has become an increasingly important field in recent years. It is used in a wide variety of applications, from chatbots to automated customer service systems. As more and more data becomes available, we can expect to see even more progress in the field of NLP in the future.

How to represent words?

Direct Way

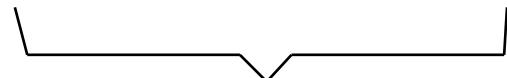
- Student
- [‘s’, ‘t’, ‘u’, ‘d’, ‘e’, ‘n’, ‘t’]: [18, 19, 20, 3, 4, 13, 19]

Problems

- Different lengths
- Not similar for similar words in meaning

One-hot coding

单词	序号	编码
A	0	[1, 0, 0, 0,, 0]
An	1	[0, 1, 0, 0,, 0]
About	2	[0, 0, 1, 0,, 0]
.....
Zero	14901	[0, 0, 0, 0,, 1]



Disadvantages:

Waste of space

Can not reflect similarity

14902

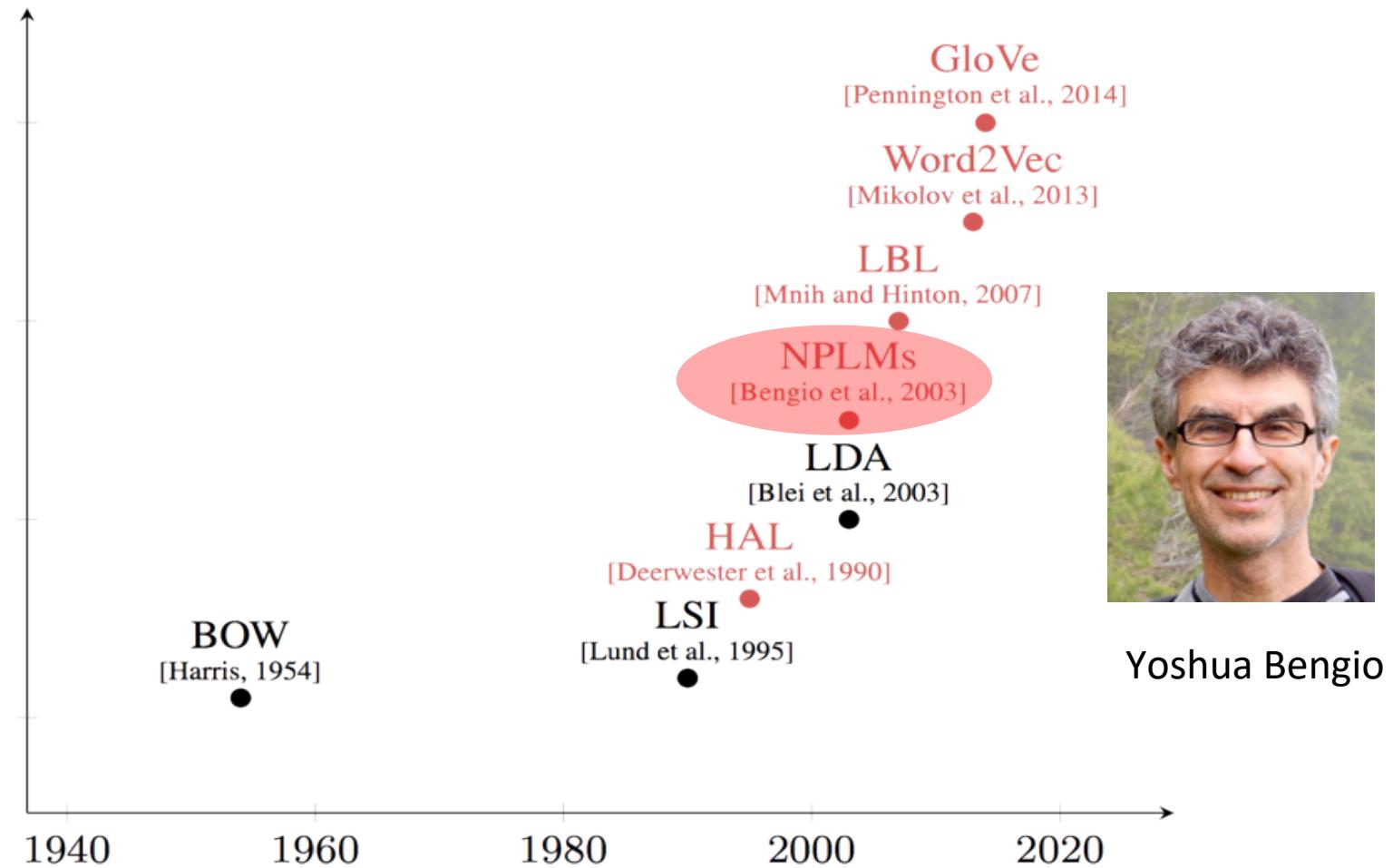
What we want

单词	词向量
star	[0.5, 0.7, -0.3, 0.2]
sun	[0.5, -0.8, 0.6, 0.4]
moon	[0.49, 0.6, -0.3, 0.1]

Features

- Dense vectors
- Similarity

Breakthroughs

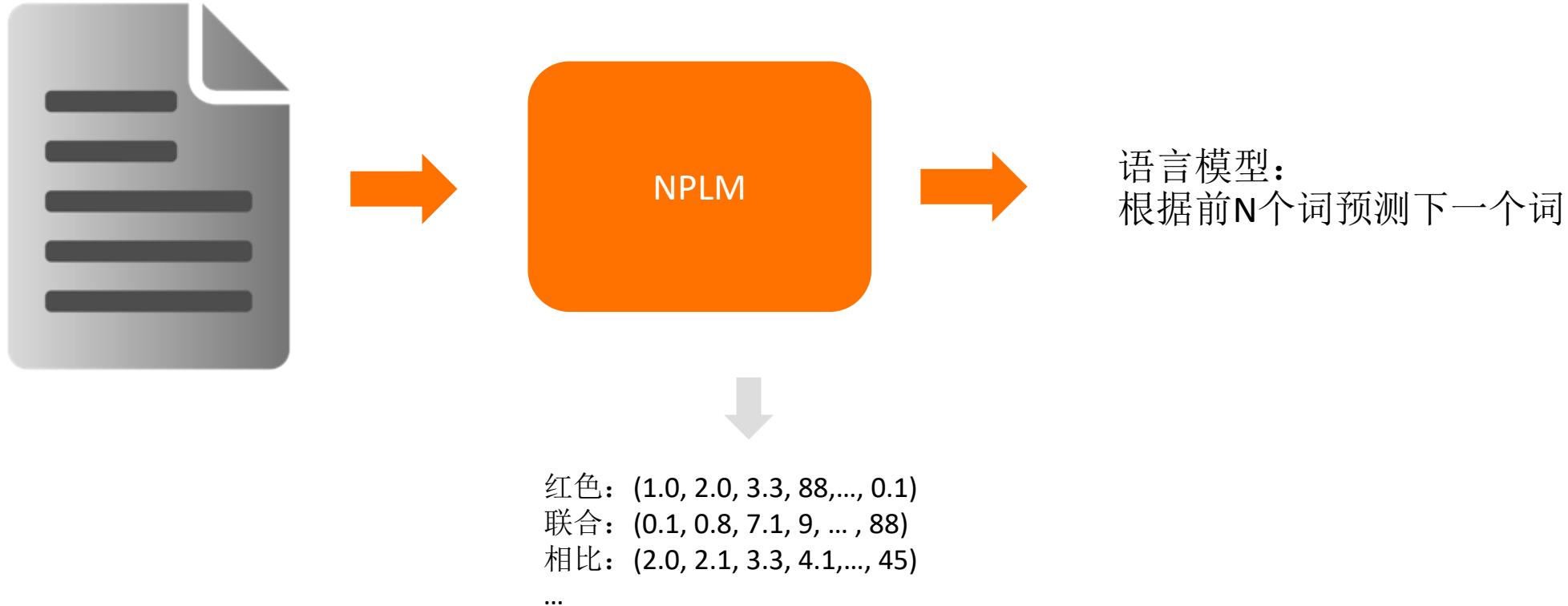


NPLM: An Unsupervised Learning Model

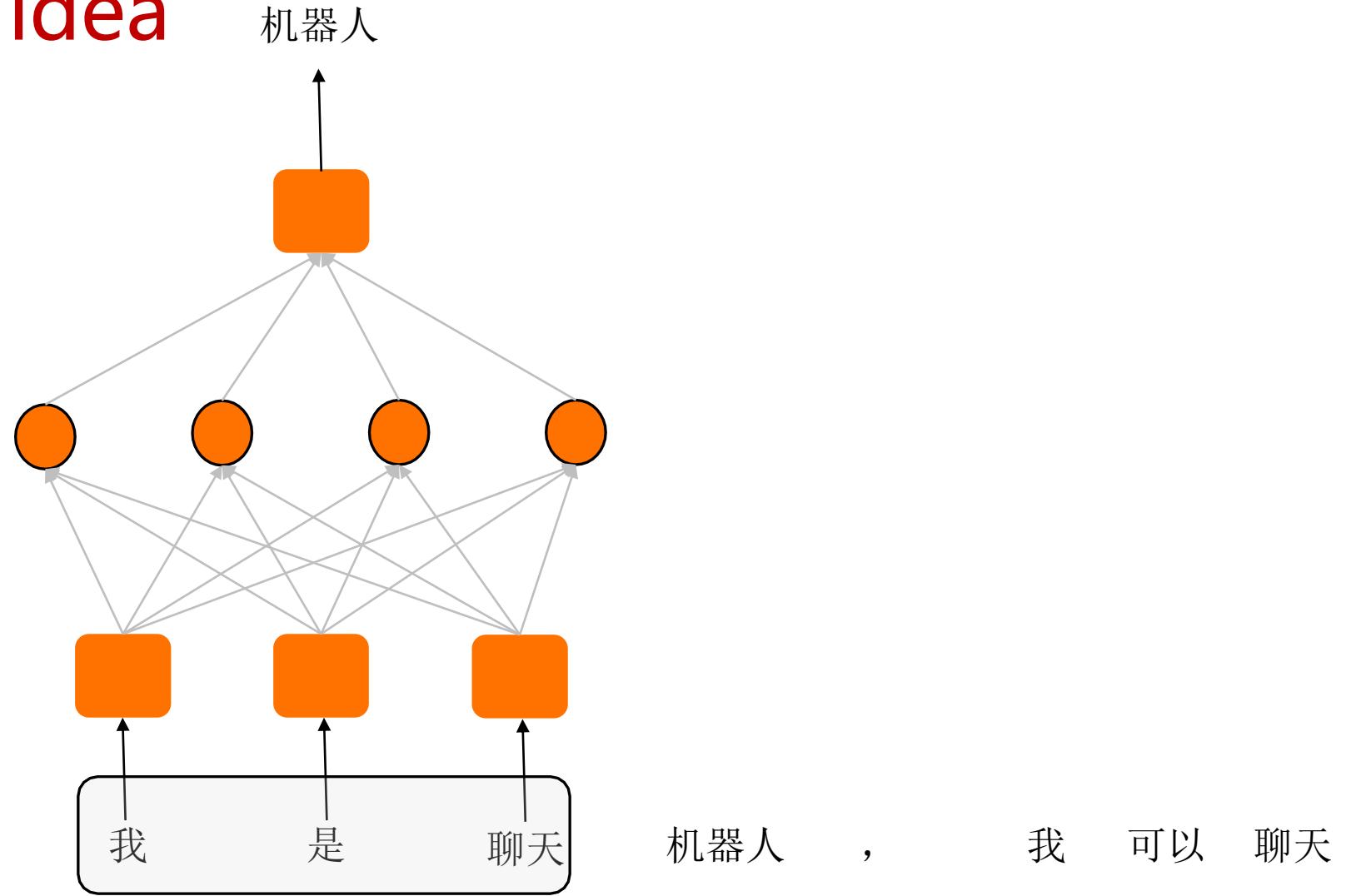


红色: (1.0, 2.0, 3.3, 88,..., 0.1)
联合: (0.1, 0.8, 7.1, 9, ..., 88)
相比: (2.0, 2.1, 3.3, 4.1,..., 45)
...
理想: (0.1, 0.8, 7.1, 9, ..., 88)

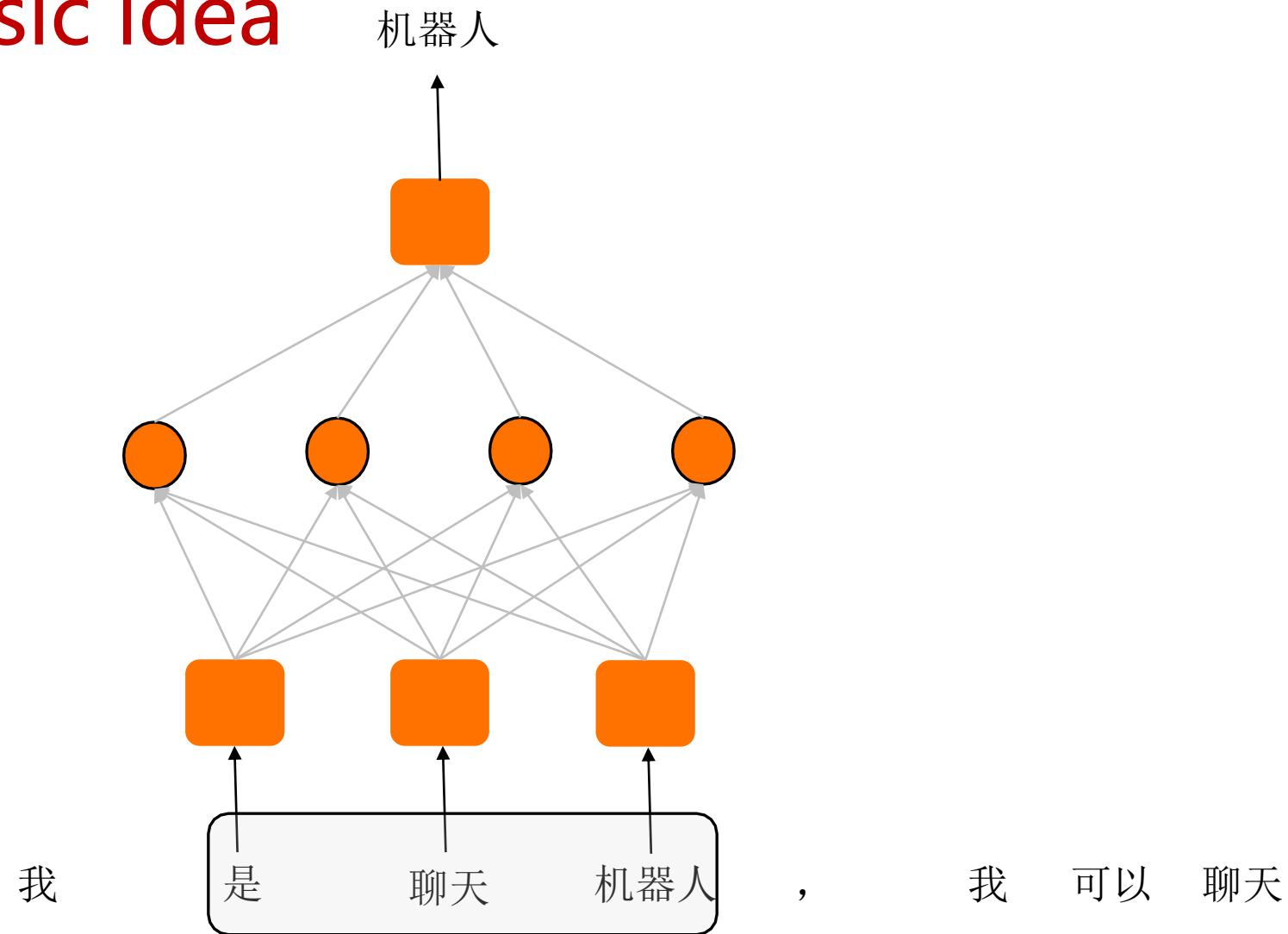
NPLM: An Unsupervised Learning Model



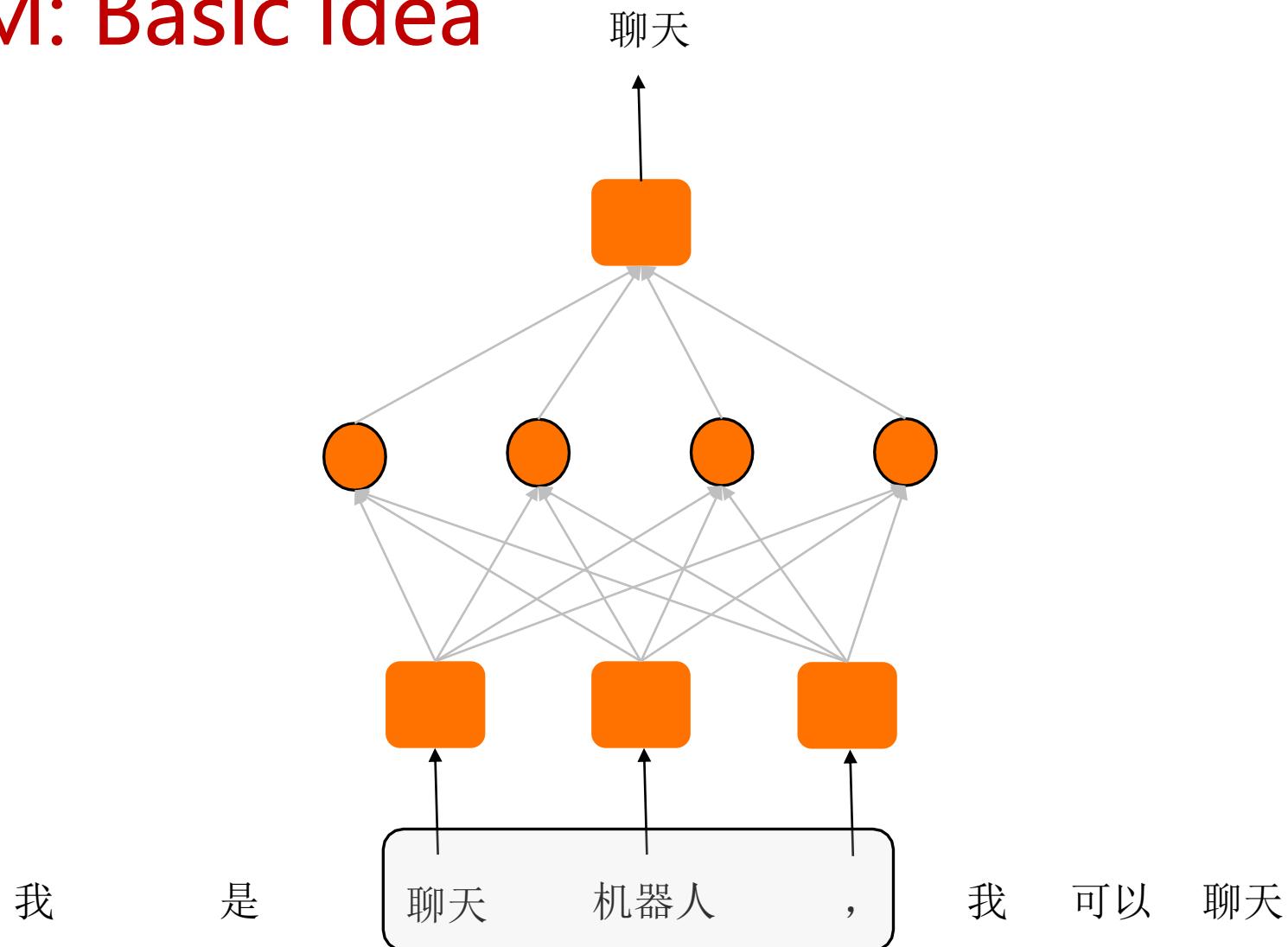
NPLM: Basic idea



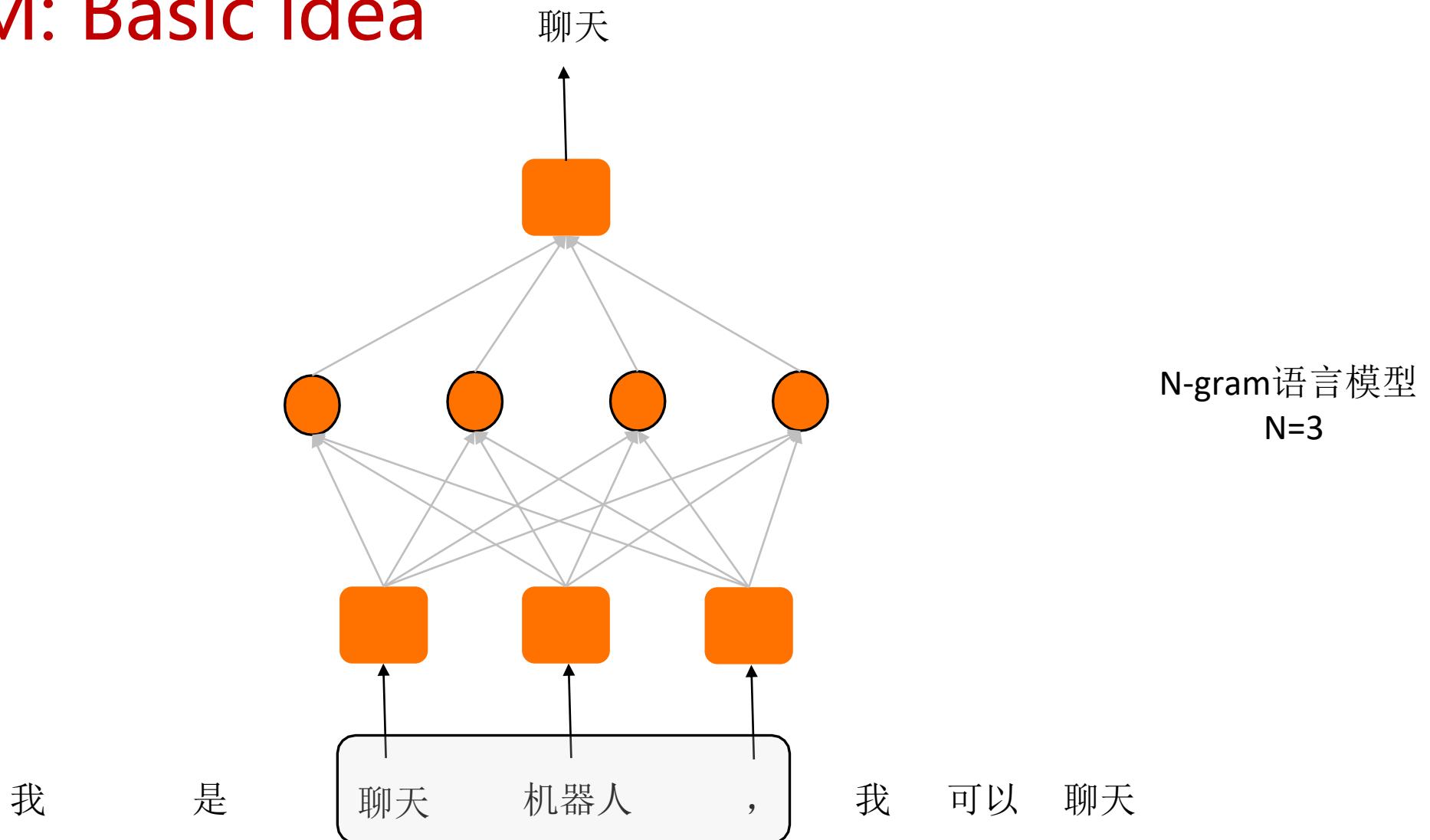
NPLM: Basic idea



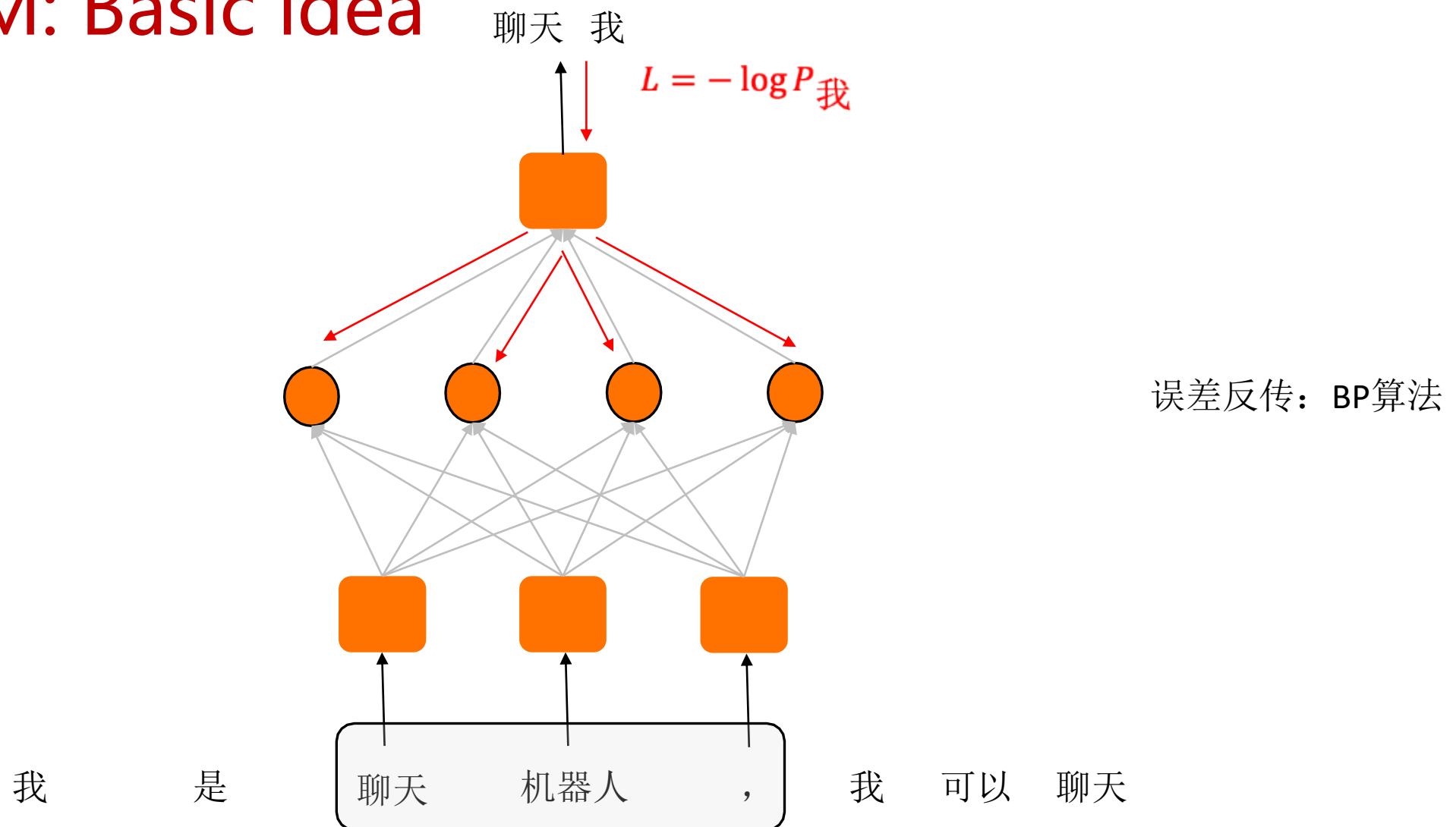
NPLM: Basic idea



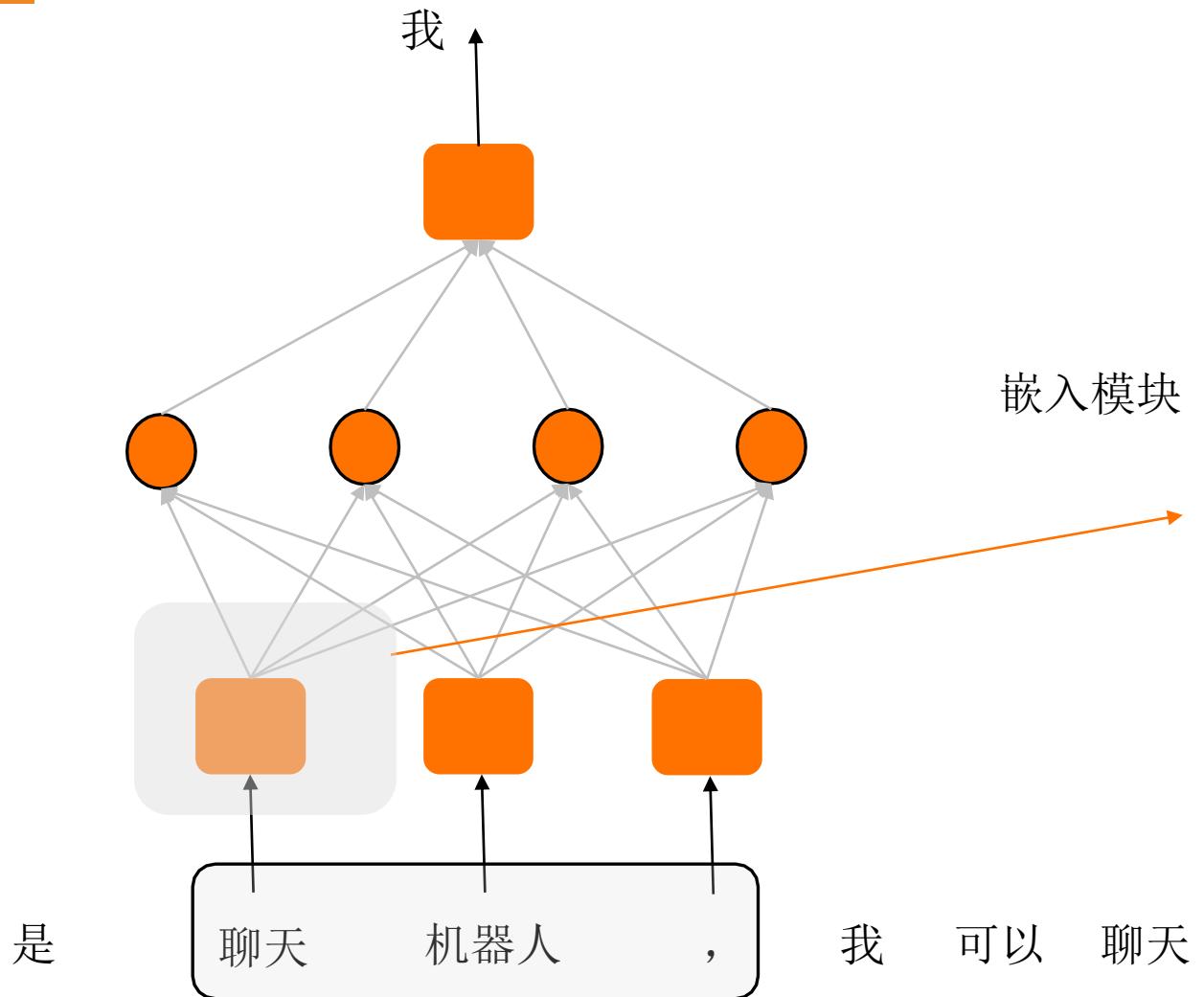
NPLM: Basic idea



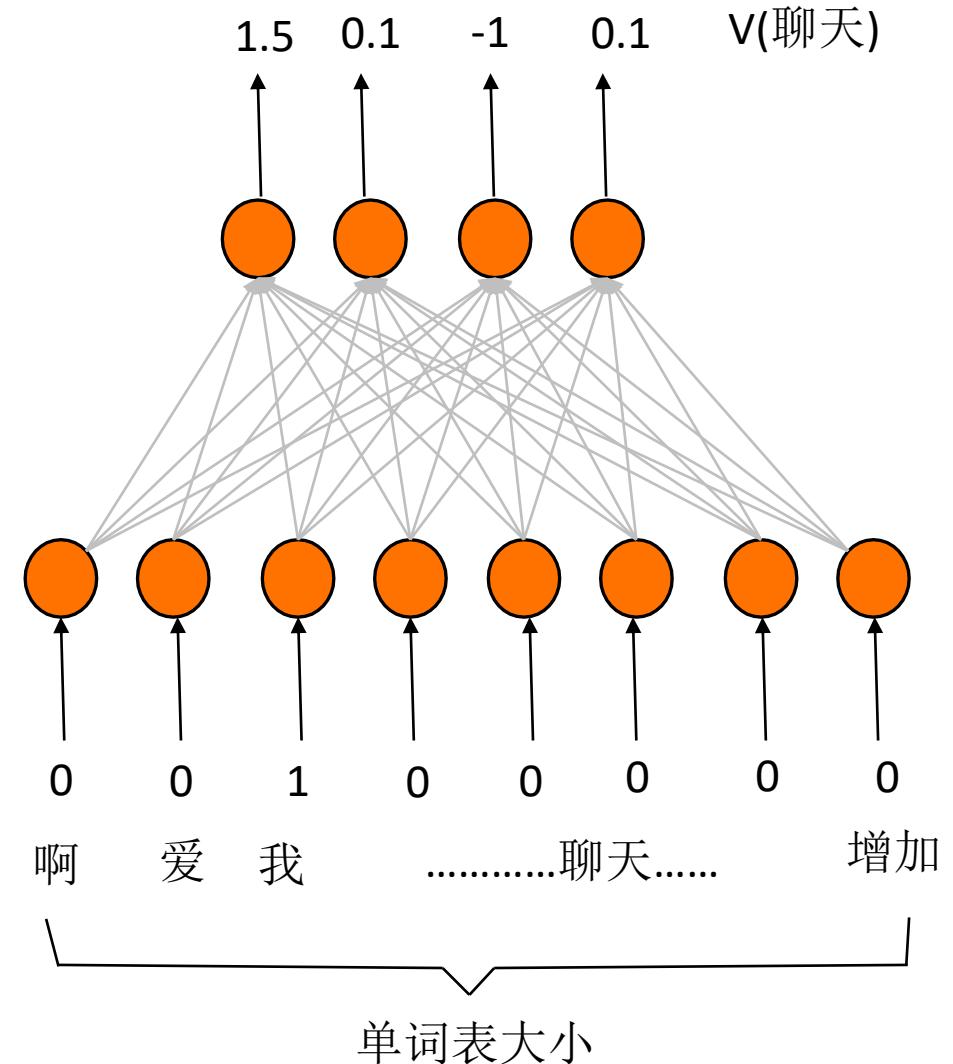
NPLM: Basic idea



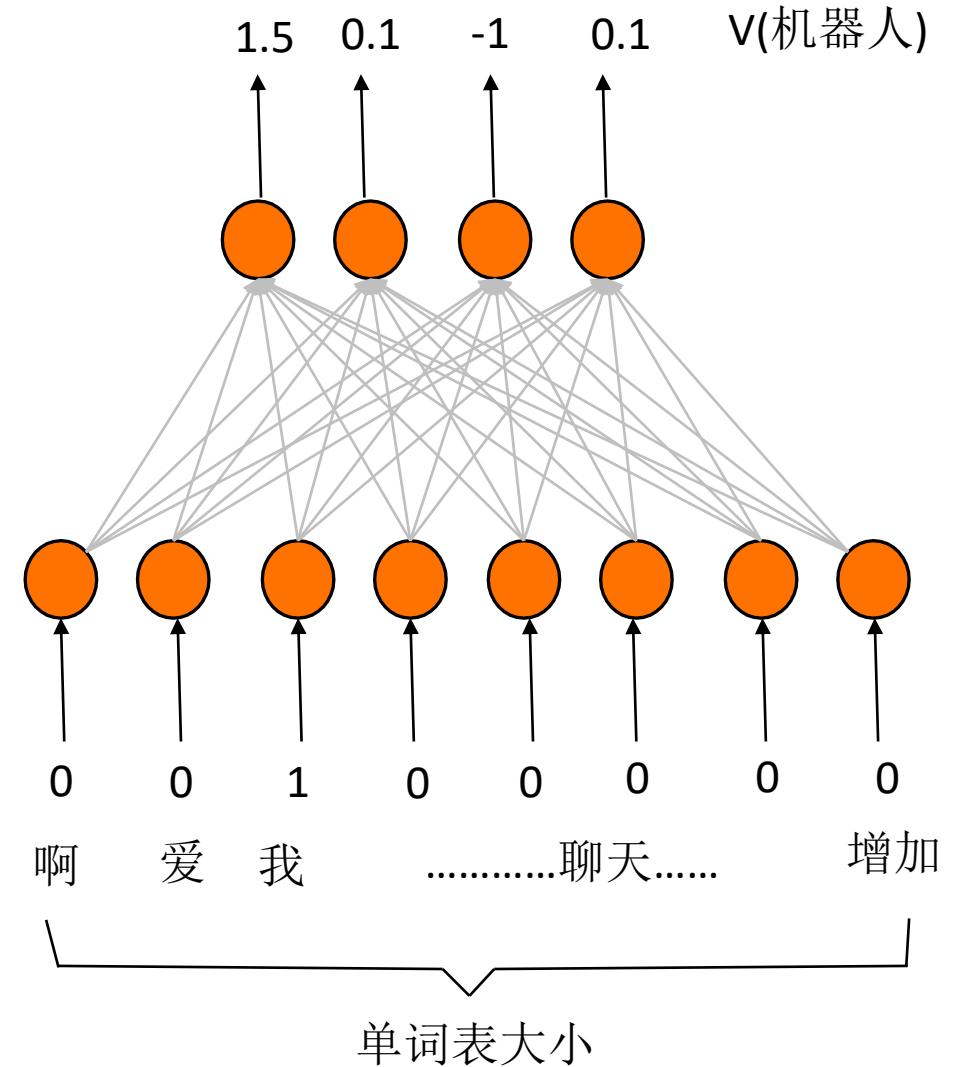
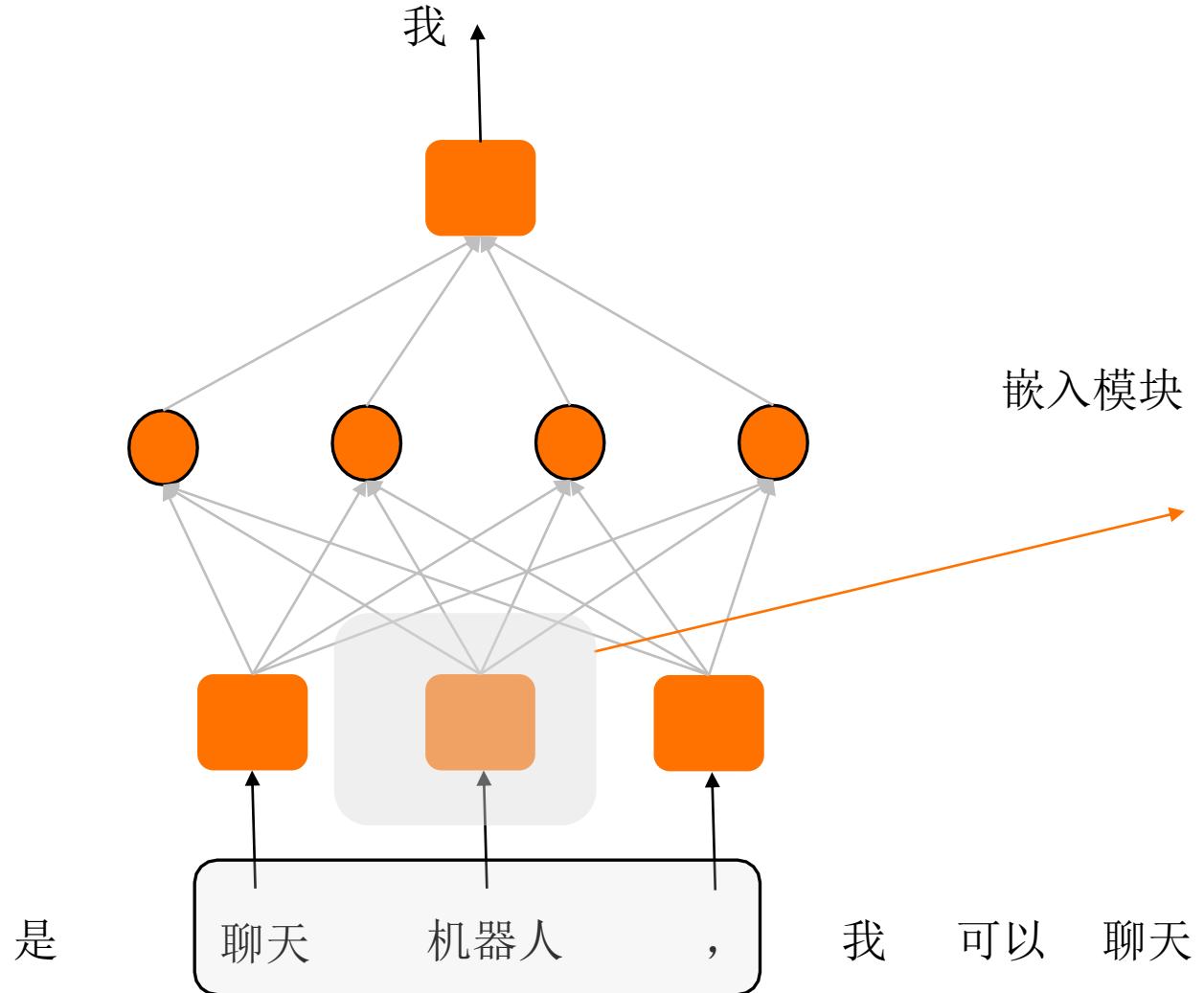
NPLM: detail



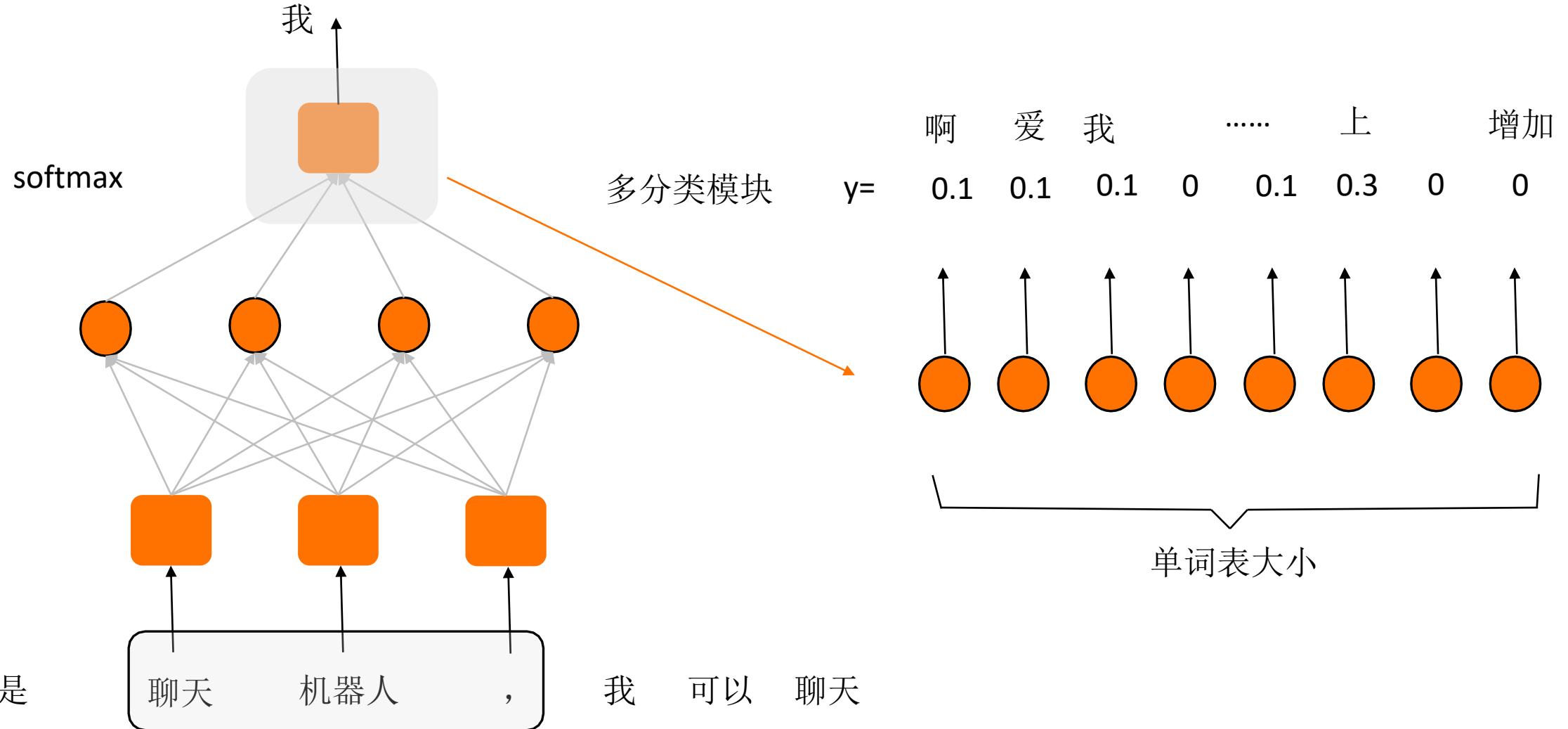
嵌入模块



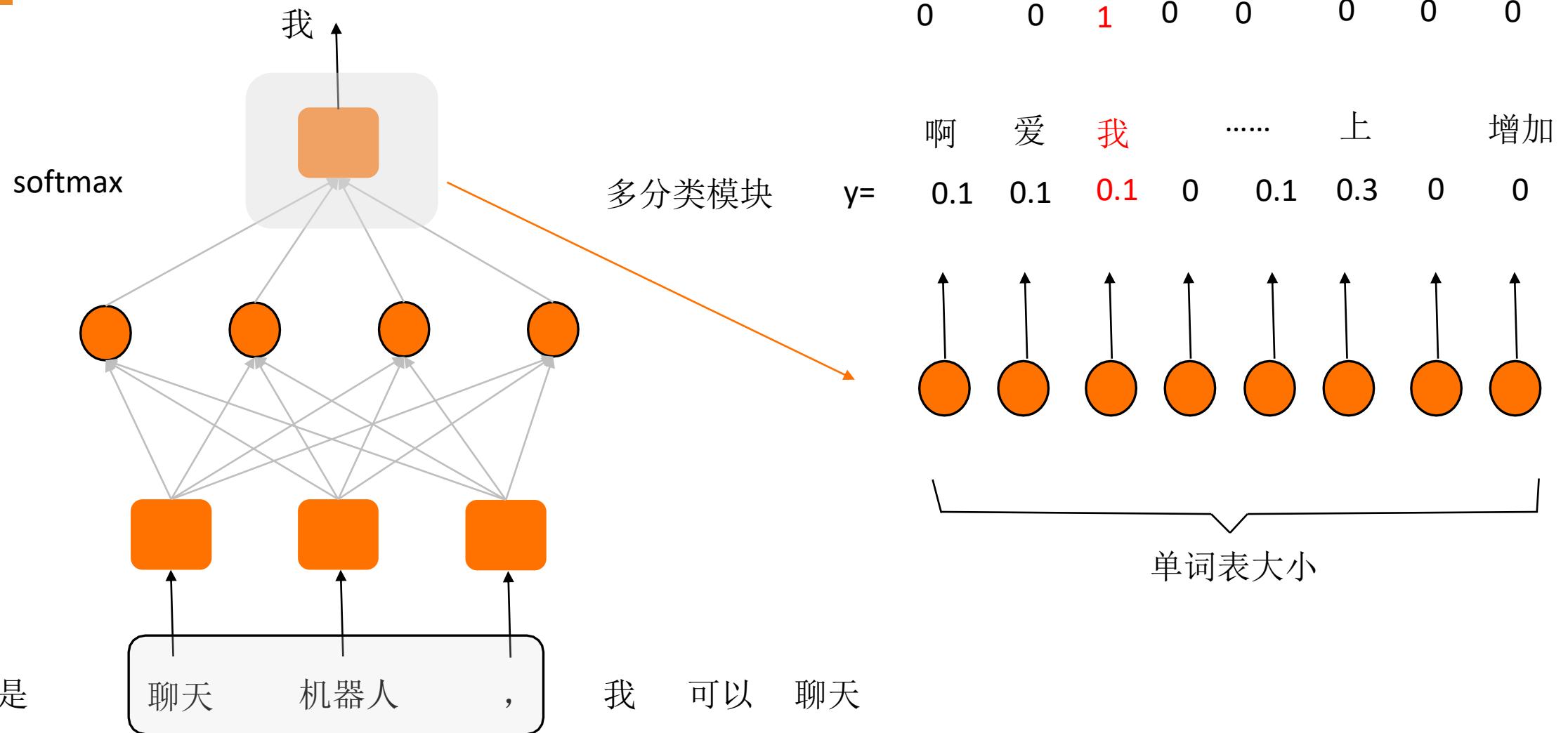
NPLM: detail



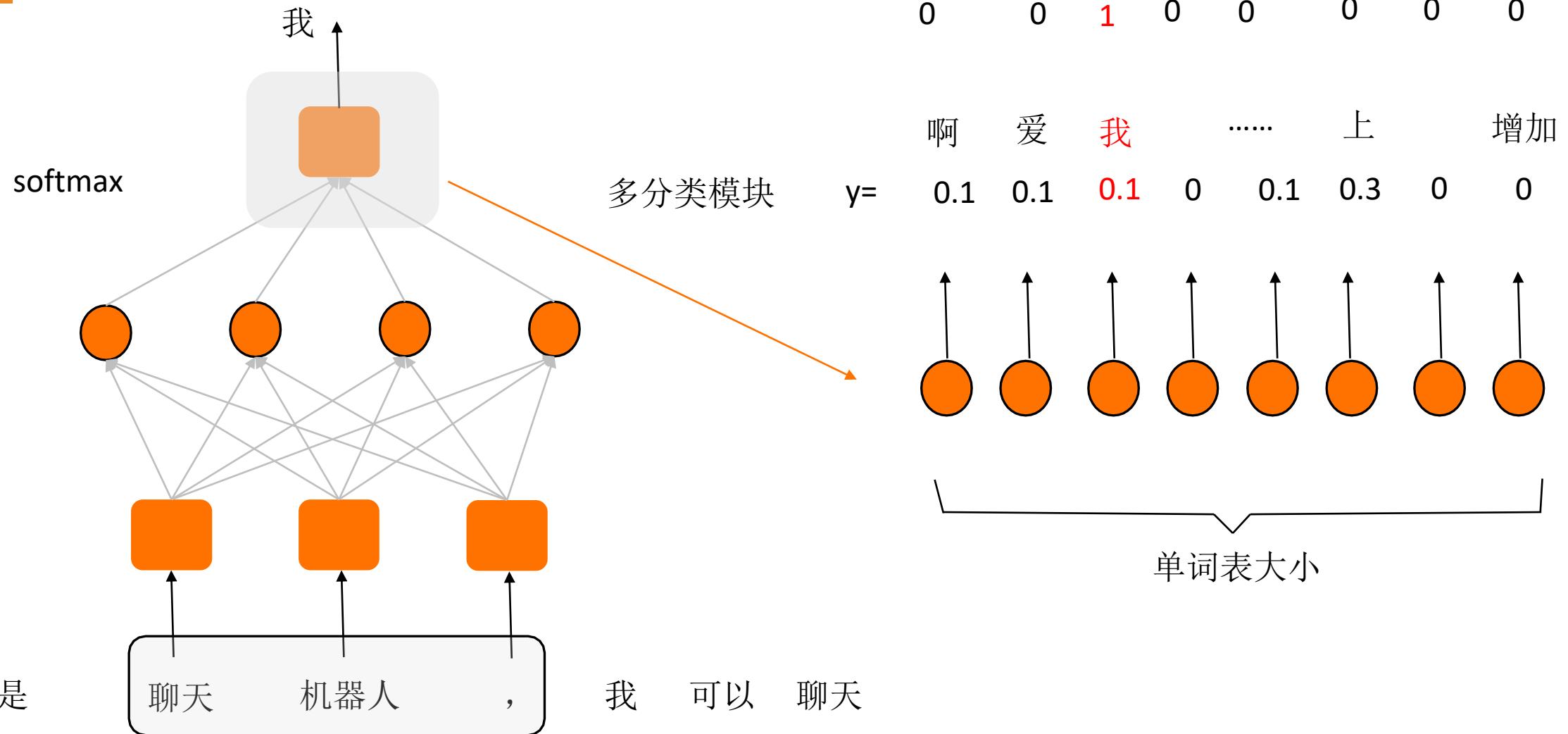
NPLM: detail



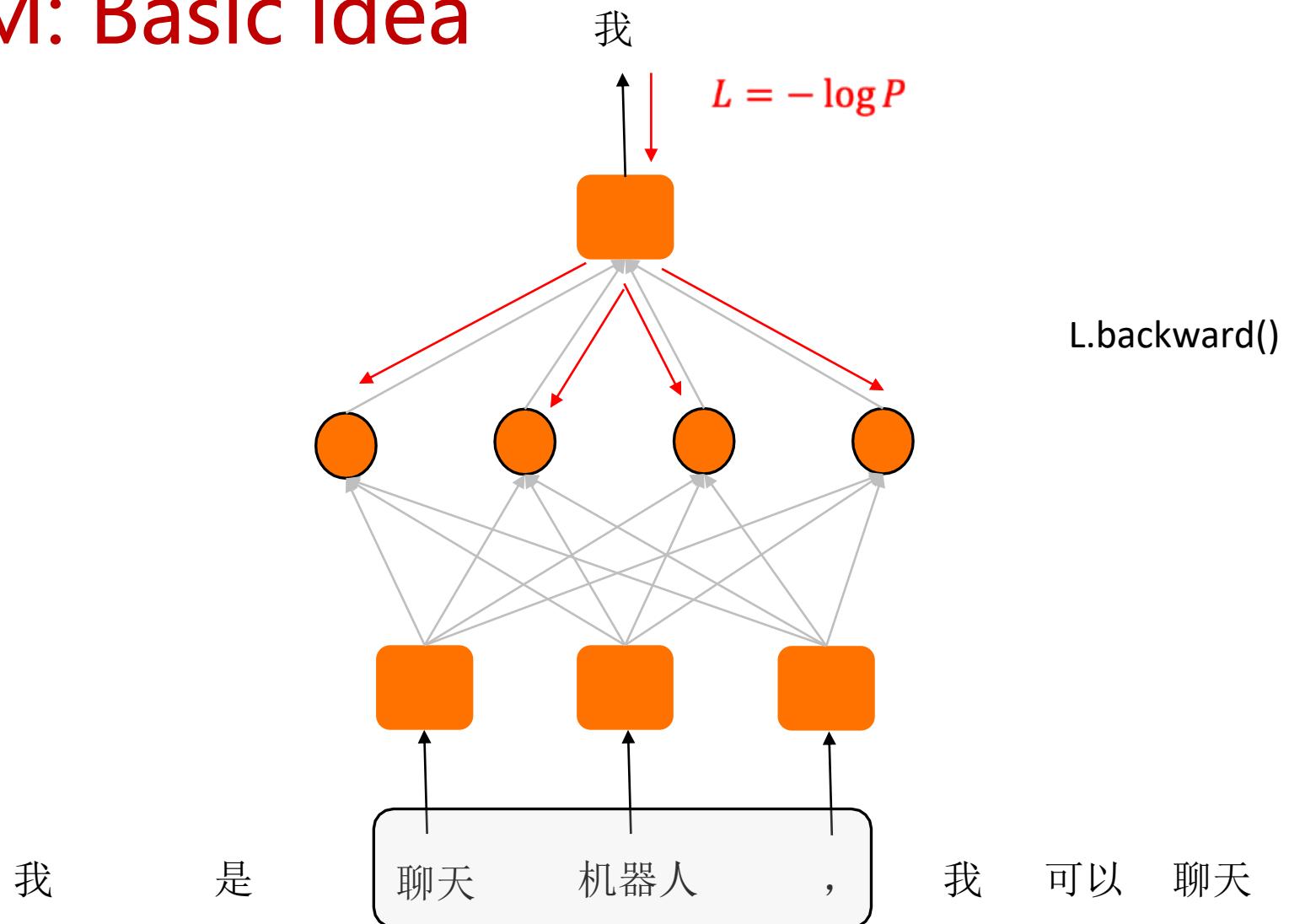
NPLM: detail



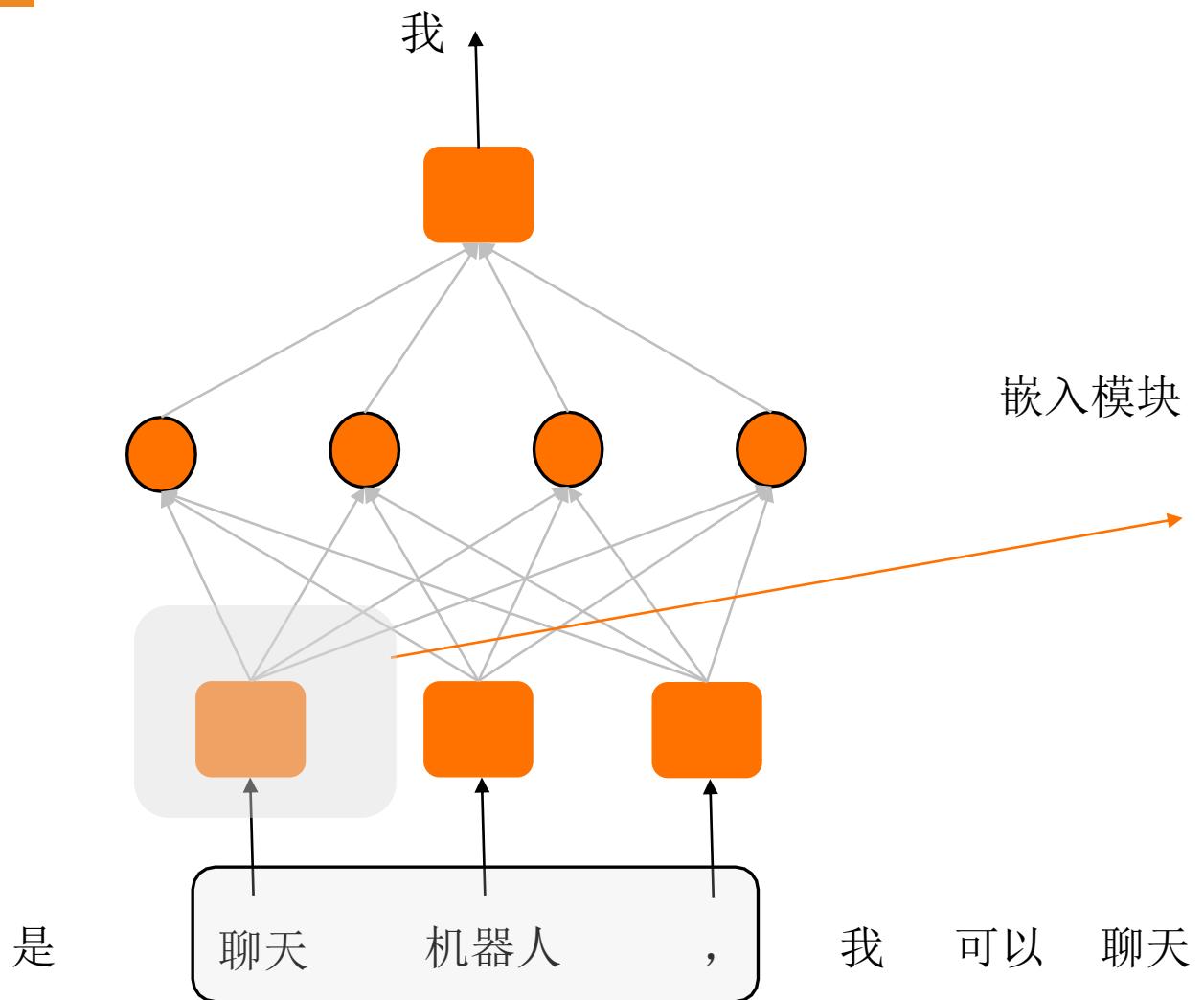
NPLM: detail



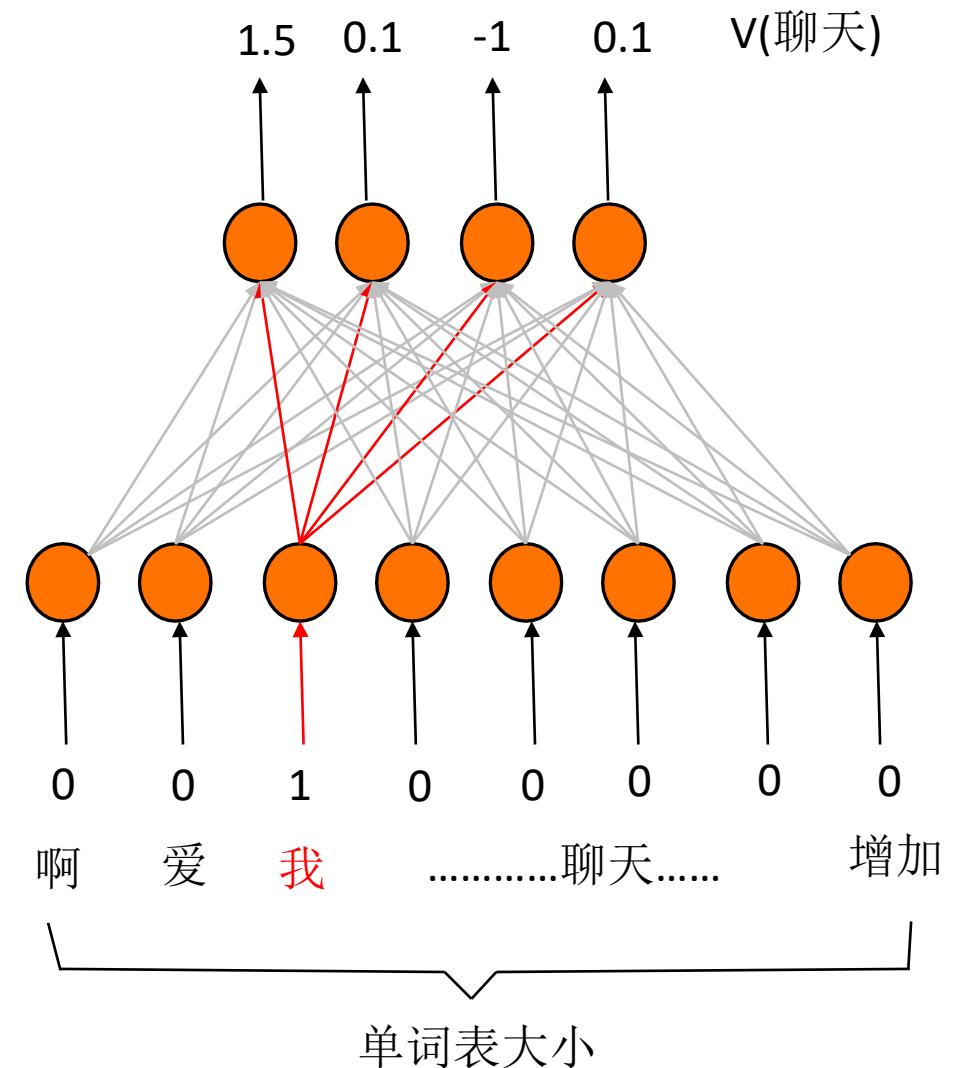
NPLM: Basic idea



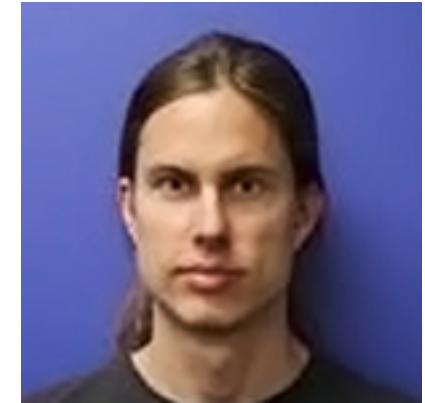
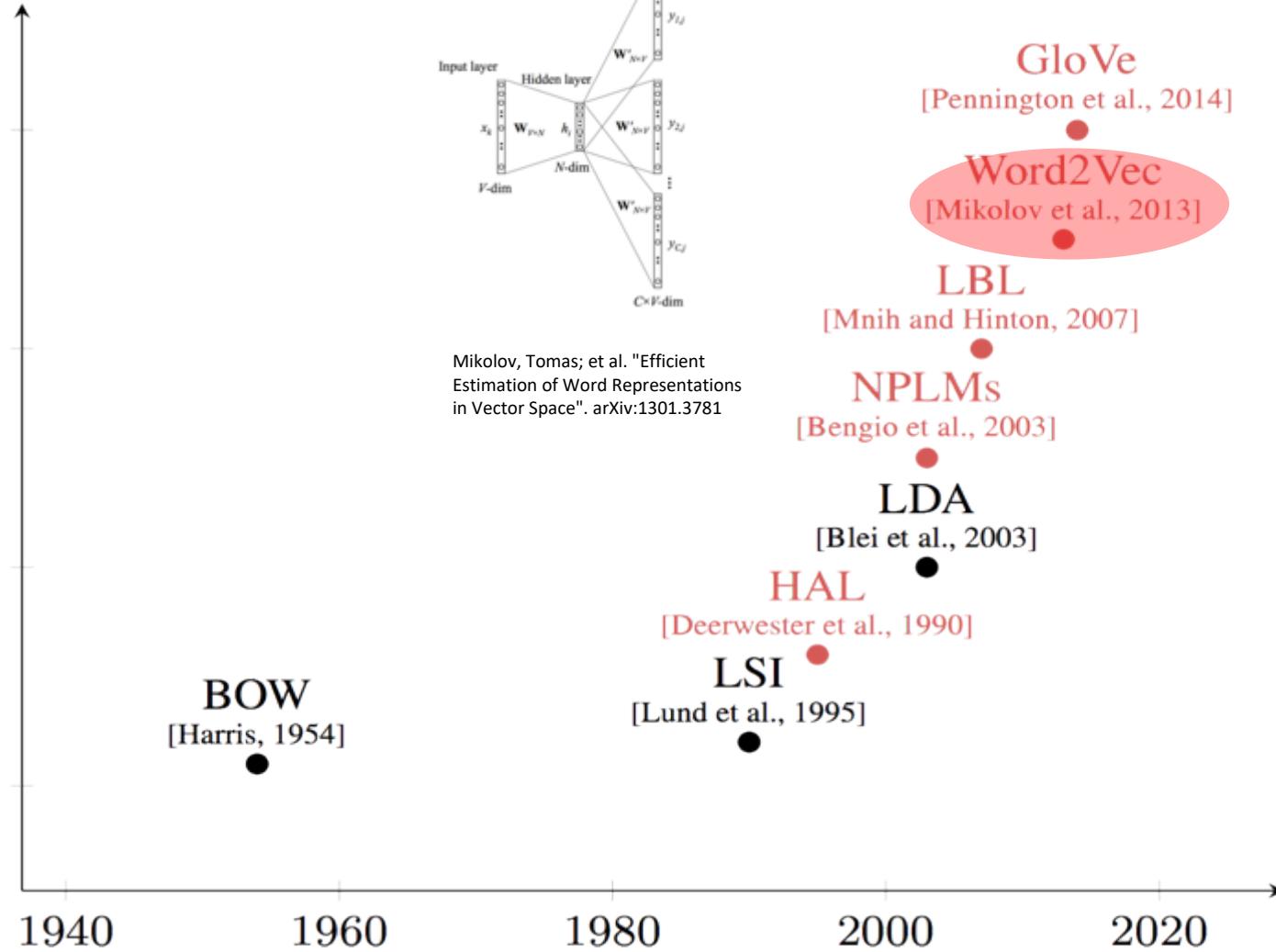
Where is the vector?



嵌入模块



Breakthroughs



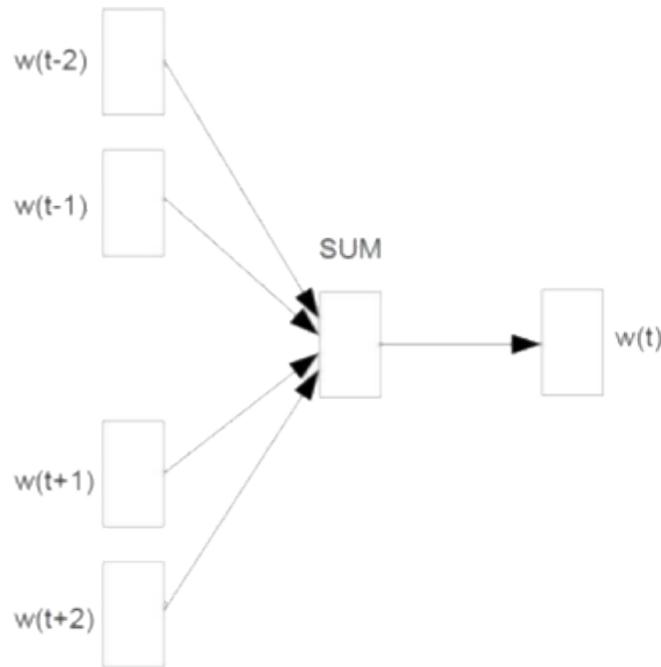
Tomas Mikolov

Word2Vec

输入

投影

输出



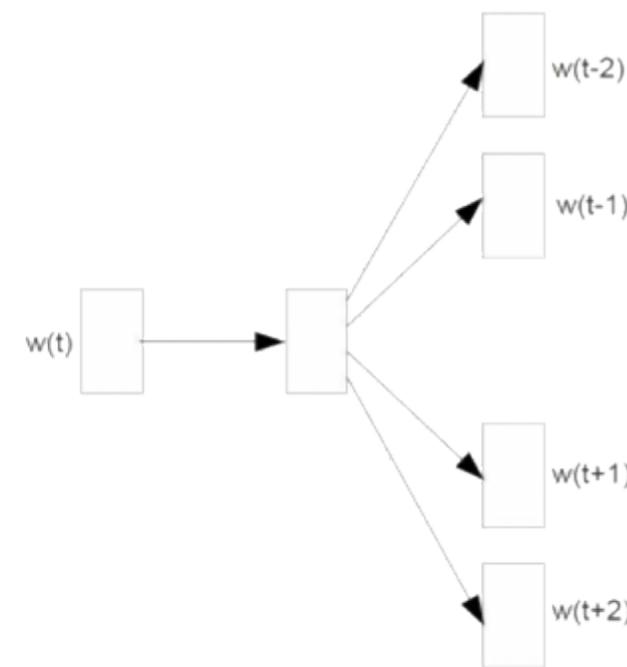
CBOW 模型

已知上下文预测当前词

输入

投影

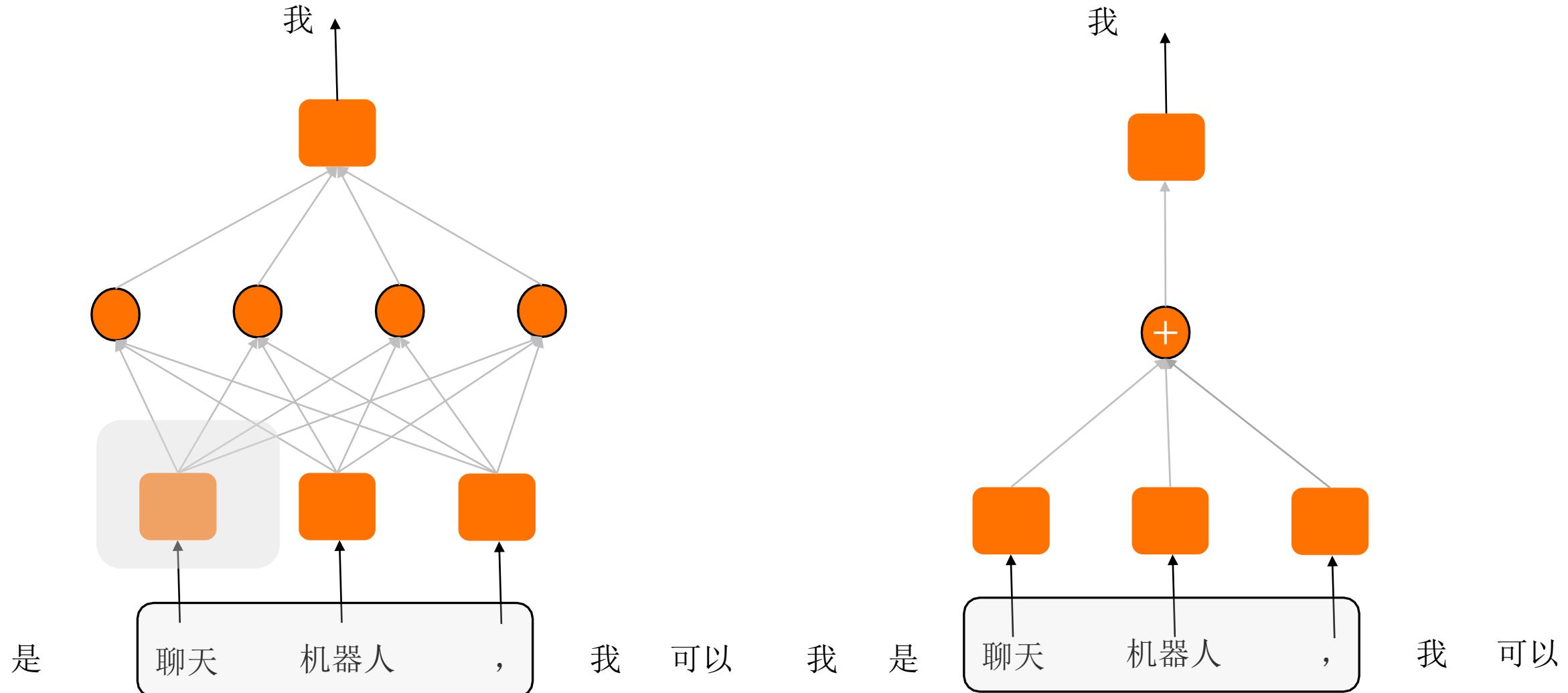
输出



Skip-gram 模型

已知当前词预测上下文

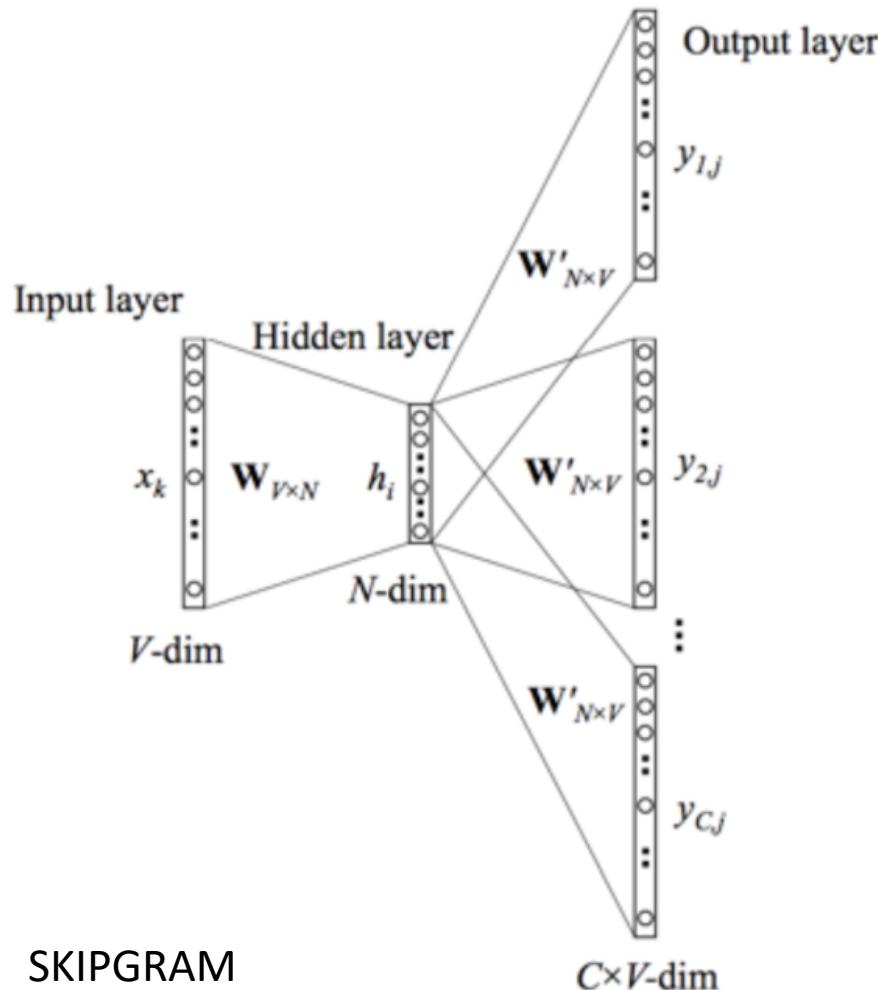
CBOW(Continuous Bag of Words)



Mikolov, Tomas; et al. "Efficient Estimation of Word Representations in Vector Space". arXiv:1301.3781



Skipgram



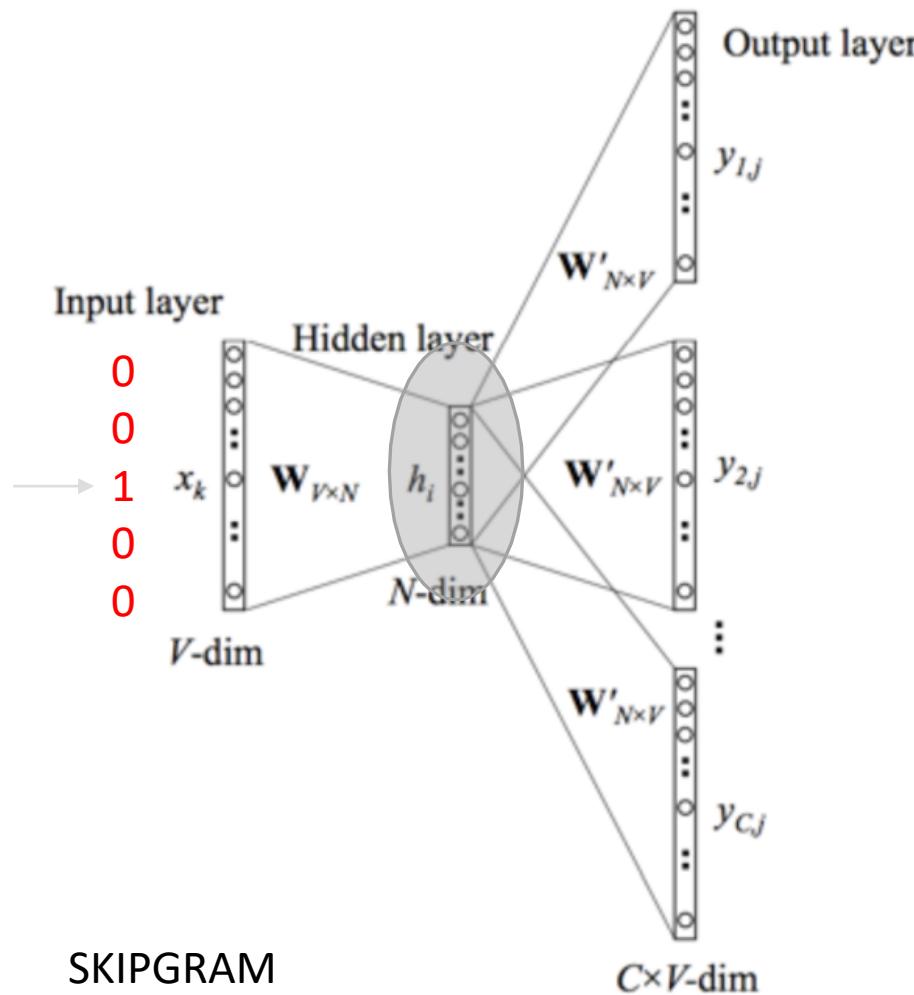
- 假设输入的句子是：

进入 21世纪 以来，网络 信息 技术 已经 渗透到 经济发展 和 社会生活 的 方方面面

- 假设窗口为3，则得到的分块如下：

进入_21世纪_以来_网络_信息_技术_已经
21世纪_以来_网络_信息_技术_已经_渗透到
以来_网络_信息_技术_已经_渗透到_经济发展
网络_信息_技术_已经_渗透到_经济发展_和
信息_技术_已经_渗透到_经济发展_和_社会生活
技术_已经_渗透到_经济发展_和_社会生活_的
已经_渗透到_经济发展_和_社会生活_的_方方面面

Skipgram



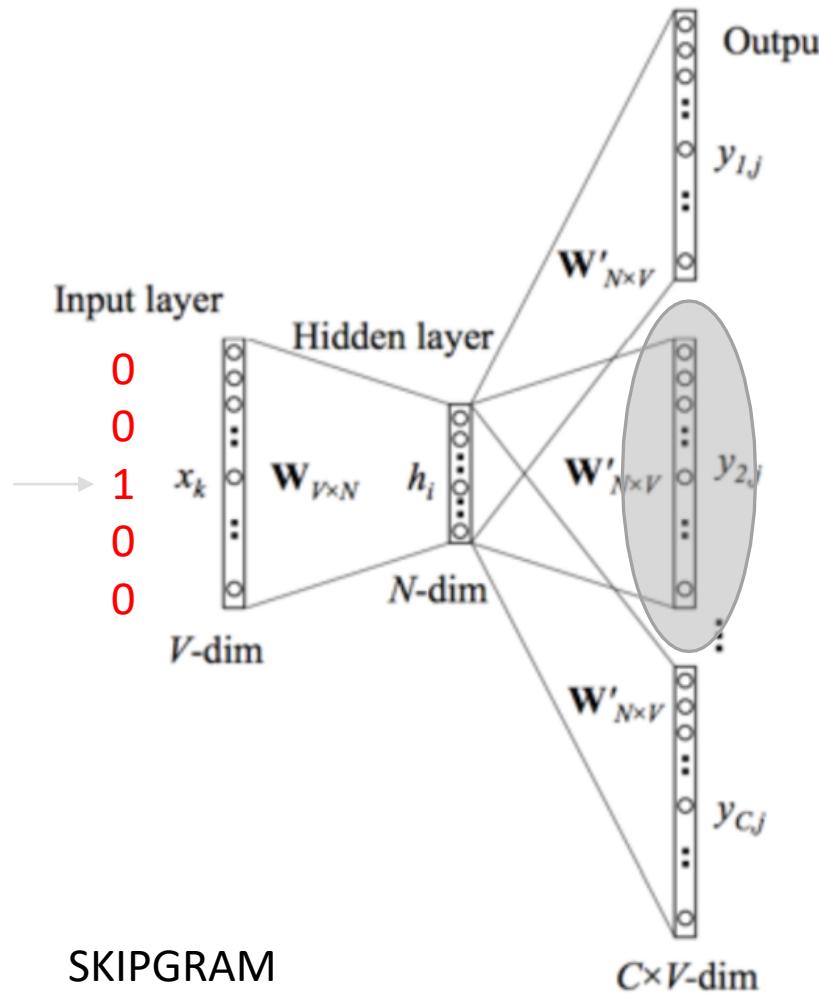
- 考虑第一个序列

进入_21世纪_以来_网络_信息_技术_已经

- 网络就是中心词，假设它的one-hot编码是： $x_k=00100$

$$h_i = W_{V \times N} x_k$$

Skipgram



- 考虑第一个序列

进入_21世纪_以来_网络_信息_技术_已经

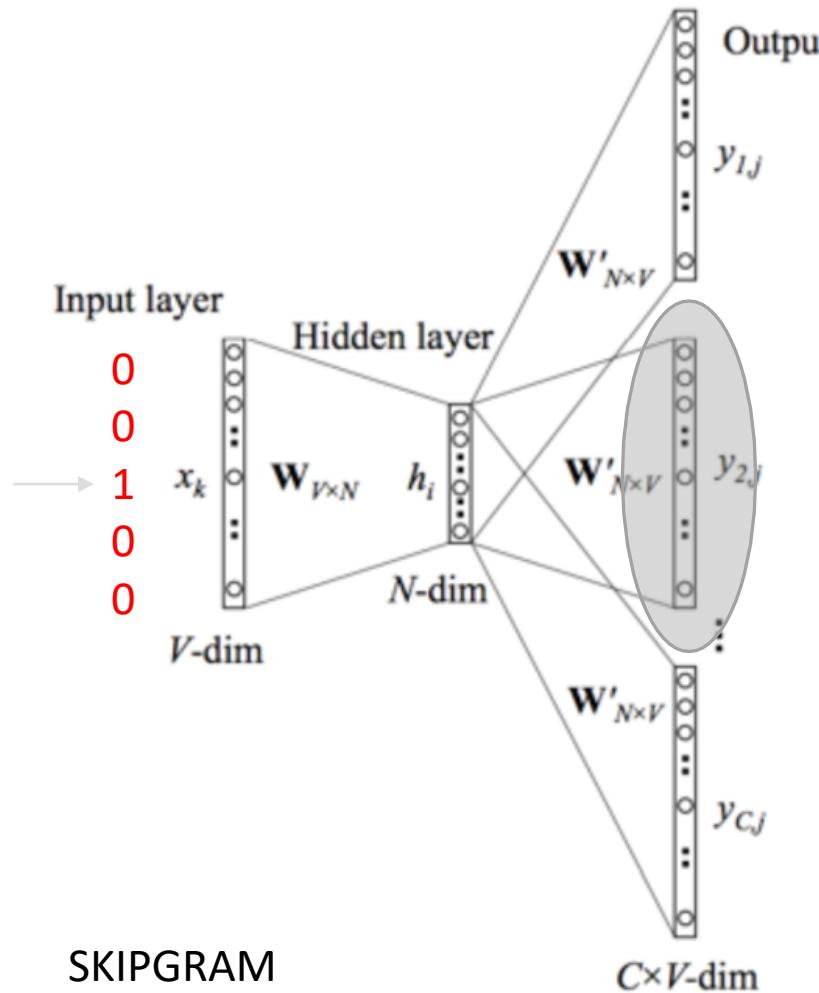
- 网络就是中心词，假设它的one-hot编码是： $x_k=00100$

$$h_i = W_{V \times N} x_k$$

$$y_{i,j} = S(W'_{N \times V} h_i) \quad S(x_i) = \frac{\exp(x_i)}{\sum_i \exp(x_i)}$$

- $y_{i,j}$ 表示的是上下文中的第*i*个单词隶属于单词表中第*j*个单词的概率

Skipgram



- 考虑第一个序列

进入_21世纪_以来_网络_信息_技术_已经

- 网络就是中心词，假设它的one-hot编码是： $x_k=00100$

$$h_i = W_{V \times N} x_k$$

$$y_{i,j} = S(W'_{N \times V} h_i) \quad S(x_i) = \frac{\exp(x_i)}{\sum_i \exp(x_i)}$$

- $y_{i,j}$ 表示的是上下文中的第*i*个单词隶属于单词表中第*j*个单词的概率
- 给定了上下文，要让这个y最大，从而梯度下降，更新 \mathbf{W}'
- \mathbf{W} 就存储了所有单词的词向量