

Natural Language Processing and Its Applications

#L3

Neural Network based Language Model

袁彩霞

yuancx@bupt.edu.cn

智能科学与技术中心

Outline

- Motivation
- Word Representation
- Neural network based language models
 - Feedforward NNLM
 - Recurrent NNLM

N-grams

- Standard approach to language modeling
- Task: compute probability of a sentence W

$$P(w_1, w_2, w_3, \dots, w_n) = \prod_i P(w_i | w_1, w_2, \dots, w_{i-1})$$

- Often simplified to trigrams:

$$P(w_1, w_2, w_3, \dots, w_n) = \prod_i P(w_i | w_{i-2}, w_{i-1})$$

N-grams: example

- $P(\text{"this is a sentence"}) = P(\text{this}) \times P(\text{is} | \text{this}) \times P(\text{a} | \text{this, is}) \times P(\text{sentence} | \text{is, a})$

- The probabilities are estimated from counts:

$$P(a | \text{this, is}) = \frac{C(\text{this is a})}{C(\text{this is})}$$

- Smoothing is used to redistribute probability to unseen events (this avoids zero probabilities)

n-gram模型的局限

- 问题1：数据稀疏问题
 - 理论上，模型阶数越高越好，但由于数据稀疏，N-gram模型中n达到一定值后，n越大性能反而越差(e.g., <6)，有没有可以算高阶的模型？
 - 同样，由于数据稀疏问题，平滑很重要，有没有不需要平滑就可以直接用的模型？

n-gram模型的局限

- 问题2：没有考虑词的含义
 - The cat is walking in the bedroom的训练样本对
A dog was running in a room的句子概率无贡献
 - 没有相同的bigram
 - 更细节： $p(\text{eat}|\text{cat})$ 和 $p(\text{eat}|\text{dog})$ 无关
 - cat 和 dog无关，所以两个概率无关，因此，在某些语料中，这两个值可能差别很大
 - 词相似，概率有理由相似： $p(\text{eat}|\text{cat}) \sim p(\text{eat}|\text{dog})$

- 基于符号的词表示方法做不到！
- 词的符号表示？
- 词的语义表示？
 - How to represent word meaning?

How do we represent the meaning of a word?

- **Definition: Meaning (Webster dictionary)**
 - the idea that is represented by a word, phrase, etc.
 - the idea that a person wants to express by using words, signs, etc.
 - the idea that is expressed in a work of writing, art, etc.
 - 词义：词的含义。(也有词典解释为：词语的意义)

How to represent meaning in a computer?

- Common answer: **Use a taxonomy** like WordNet that has hypernyms (is-a) relationships and synonym sets

[Synset('procyonid.n.01'),
Synset('carnivore.n.01'),
Synset('placental.n.01'),
Synset('mammal.n.01'),
Synset('vertebrate.n.01'),
Synset('chordate.n.01'),
Synset('animal.n.01'),
Synset('organism.n.01'),
Synset('living_thing.n.01'),
Synset('whole.n.02'),
Synset('object.n.01'),
Synset('physical_entity.n.01'),
Synset('entity.n.01')]

S: (adj) full, good
S: (adj) estimable, good, honorable, respectable
S: (adj) beneficial, good
S: (adj) good, just, upright
S: (adj) adept, expert, good, practiced, proficient, skillful
S: (adj) dear, good, near
S: (adj) good, right, ripe
...
S: (adv) well, good
S: (adv) thoroughly, soundly, good
S: (n) good, goodness
S: (n) commodity, trade good, good

Problems with this discrete representation

- Great as resource but missing nuances, e.g.
 - **synonyms:** adept, expert, good, practiced, proficient, skillful?
- Missing new words (impossible to keep up to date): wicked, badass, crack, ace, wizard, genius, ninja
- Subjective
- Requires human labor to create and adapt
- Hard to compute accurate word similarity

Problems with this discrete representation

- The vast majority of rule-based **and statistical NLP work regards**
- words as atomic symbols: hotel, conference, walk
- In vector space terms, this is a vector with one 1 and a lot of zeroes

[0 0 0 0 0 0 0 0 0 0 1 0 0 0 0]

- Dimensionality:
20K (speech) – 50K (PTB) – 500K (big vocabulary) –
13M (Google 1T)
- We call this a “one-hot” representation. Its problem:
motel [0 0 0 0 0 0 0 0 0 0 1 0 0 0 0] AND
hotel [0 0 0 0 0 0 0 1 0 0 0 0 0 0 0] = 0

Distributional similarity based representations

- You can get a lot of value by representing a word by means of its neighbors

“You shall know a word by the company it keeps”

(J. R. Firth 1957: 11)

- One of the most successful ideas of modern statistical NLP

government debt problems turning into *banking* crises as has happened in Europe needs unified *banking* regulation to replace the hodgepodge

- These words will represent *banking*

How to make neighbors represent words?

- Answer: With a cooccurrence matrix X
- 2 options: full document vs windows
 - Word - document cooccurrence matrix: give general topics (all sports terms will have similar entries) leading to “Latent Semantic Analysis”
 - Window around each word: captures both syntactic and semantic information

Window based cooccurrence matrix

- Window length l (e.g., 1, more common: 5 - 10)
- Symmetric (irrelevant whether left or right context)
- Example corpus:
 - I like deep learning.
 - I like NLP.
 - I enjoy flying.

Window based cooccurrence matrix

- Example corpus:
 - I like deep learning.
 - I like NLP.
 - I enjoy flying.

counts	I	like	enjoy	deep	learning	NLP	flying	.
I	0	2	1	0	0	0	0	0
like	2	0	0	1	0	1	0	0
enjoy	1	0	0	0	0	0	1	0
deep	0	1	0	0	1	0	0	0
learning	0	0	0	1	0	0	0	1
NLP	0	1	0	0	0	0	0	1
flying	0	0	1	0	0	0	0	1
.	0	0	0	0	1	1	1	0

Problems with simple cooccurrence vectors

- Increase in size with vocabulary
- Very high dimensional: require a lot of storage
- Sparsity issues → Models are less robust

Solution: Low dimensional vectors

- Idea: store “most” of the important information in a fixed, small number of dimensions: a dense vector
- Usually around 25 – 1000 dimensions
- How to reduce the dimensionality?

Method 1: Dimensionality Reduction on X

- Singular Value Decomposition of cooccurrence matrix X

$$\begin{array}{ccccc}
 \begin{array}{c} m \\ \boxed{} \\ n \\ X \end{array} & = & \begin{array}{c} r \\ \boxed{\begin{array}{c} | \\ U_1 \\ | \\ U_2 \\ | \\ U_3 \\ \vdots \end{array}} \\ n \\ U \end{array} & \begin{array}{c} r \\ \boxed{\begin{array}{ccc} S_1 & & \\ & S_2 & \\ & & S_3 \\ & 0 & \ddots \\ & & & S_r \end{array}} \\ r \\ S \end{array} & \begin{array}{c} m \\ \boxed{\begin{array}{c} \text{---} V_1 \text{---} \\ \text{---} V_2 \text{---} \\ \text{---} V_3 \text{---} \\ \vdots \end{array}} \\ r \\ V^T \end{array} \\
 \\
 \begin{array}{c} m \\ \boxed{\phantom{\hat{X}}} \\ n \\ \hat{X} \end{array} & = & \begin{array}{c} k \\ \boxed{\begin{array}{c} | \\ U_1 \\ | \\ U_2 \\ | \\ U_3 \\ \vdots \end{array}} \\ n \\ \hat{U} \end{array} & \begin{array}{c} k \\ \boxed{\begin{array}{ccc} S_1 & & \\ & S_2 & \\ & & S_3 \\ & 0 & \ddots \\ & & & S_k \end{array}} \\ k \\ \hat{S} \end{array} & \begin{array}{c} m \\ \boxed{\begin{array}{c} \text{---} V_1 \text{---} \\ \text{---} V_2 \text{---} \\ \text{---} V_3 \text{---} \\ \vdots \end{array}} \\ k \\ \hat{V}^T \end{array}
 \end{array}$$

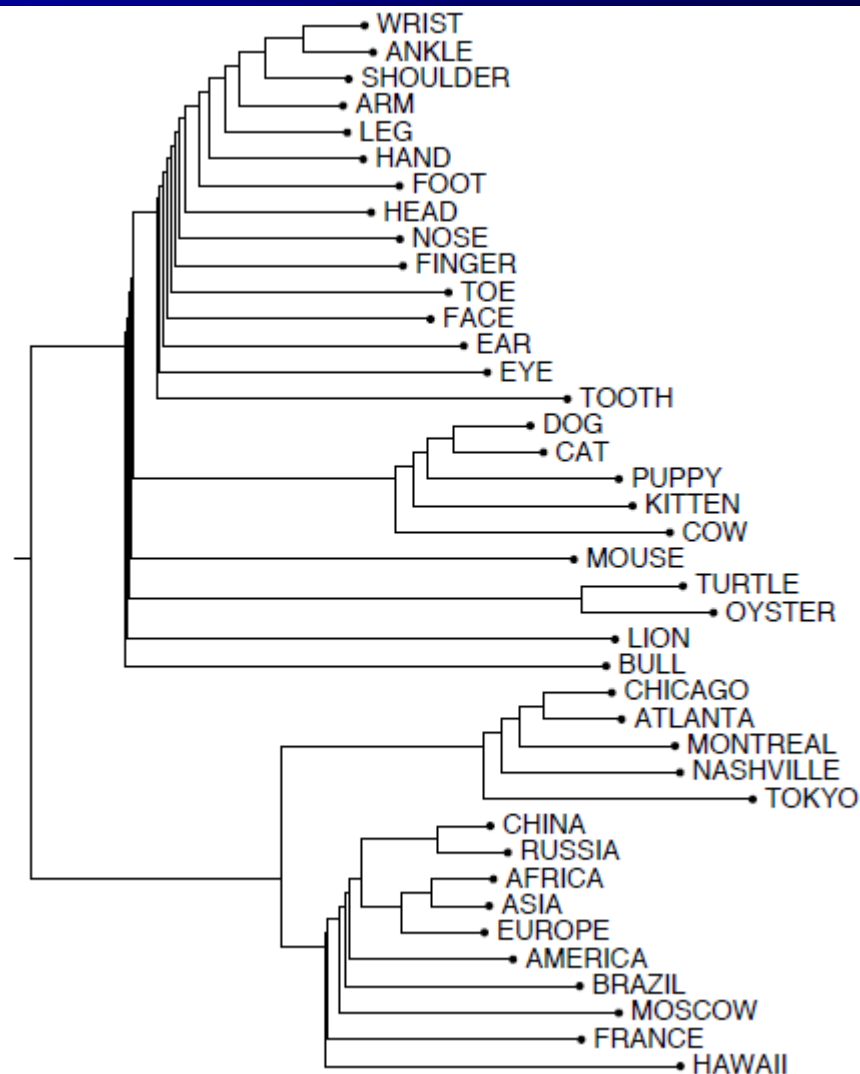
- \hat{X} is the best rank k approximation to X , in terms of least squares.

Word meaning is defined in terms of vectors

- A word is represented as a dense vector

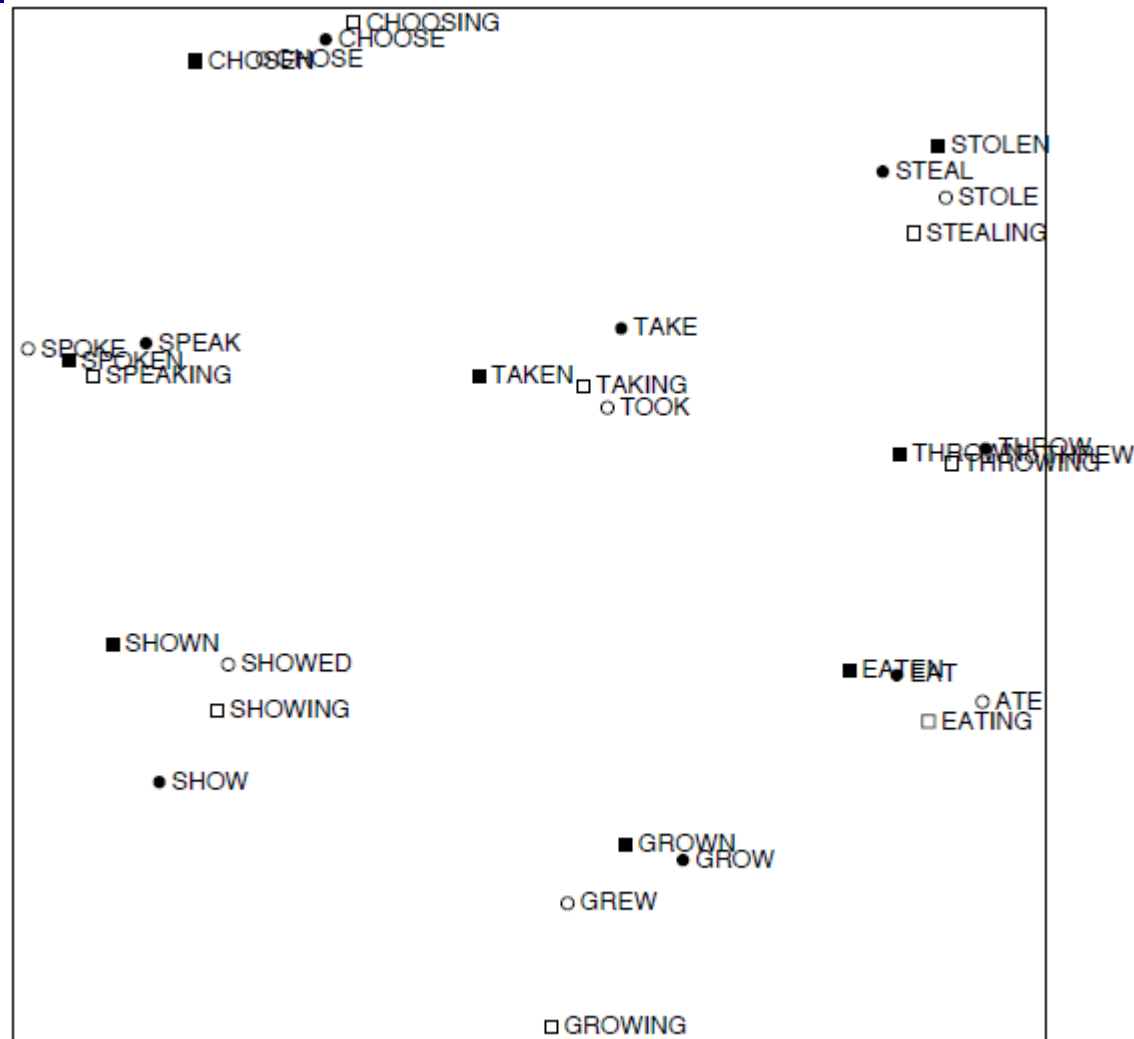
$$\textit{linguistics} = \begin{bmatrix} 0.286 \\ 0.792 \\ -0.177 \\ -0.107 \\ 0.109 \\ -0.542 \\ 0.349 \\ 0.271 \end{bmatrix}$$

Interesting semantic patterns emerge in the vectors



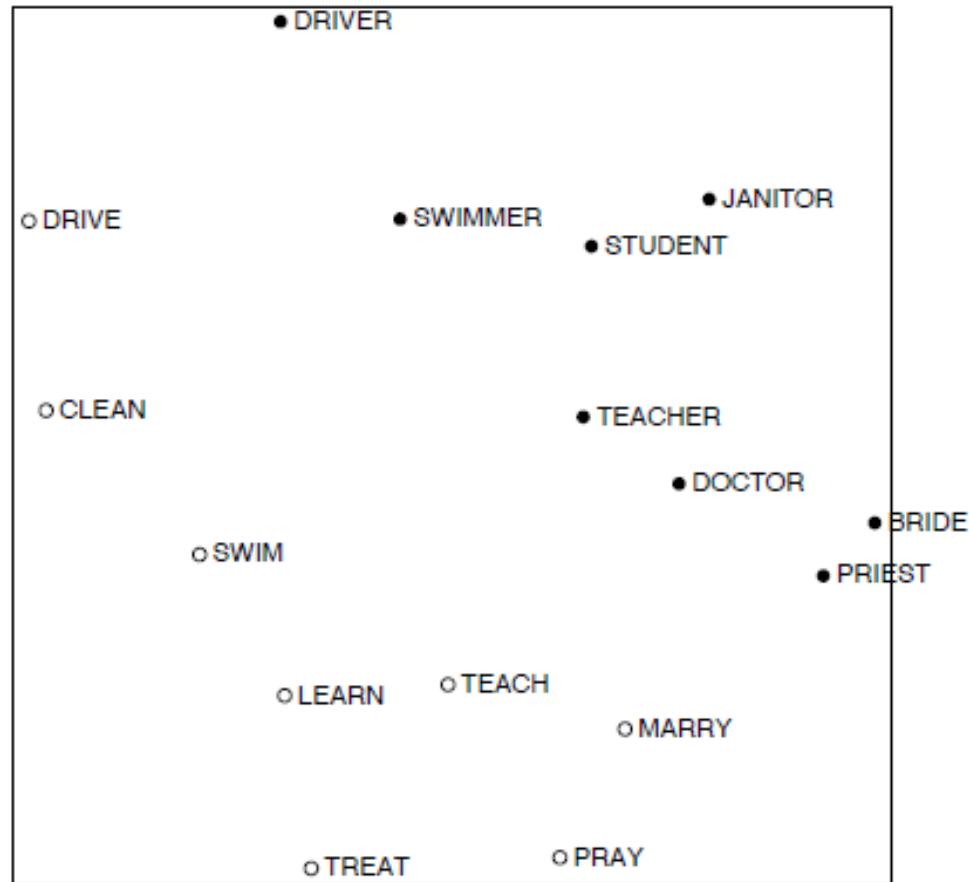
Rohde et al. An Improved Model of Semantic Similarity Based on Lexical Co-Occurrence. 2005.

Interesting syntactic patterns emerge in the vectors



Rohde et al. An Improved Model of Semantic Similarity Based on Lexical Co-Occurrence. 2005.

Interesting semantic patterns emerge in the vectors



Rohde et al. An Improved Model of Semantic Similarity Based on Lexical Co-Occurrence. 2005.

Problems with SVD

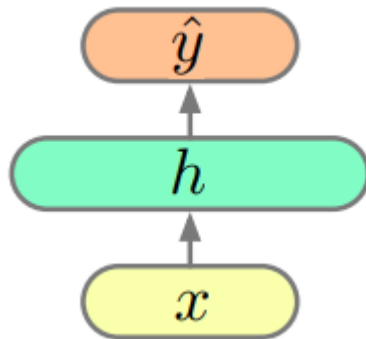
- Computational cost scales quadratically for $n \times m$ matrix: $O(mn^2)$ flops (when $n < m$)
- → Bad for millions of words or documents
- Hard to incorporate new words or documents
- Different learning regime than other deep learning models

Idea: Directly learn low-dimensional word vectors

- Learning representations by back-propagating errors. (Rumelhart et al., 1986)
- A neural probabilistic language model (Bengio et al., 2003)
- NLP (almost) from Scratch (Collobert & Weston, 2008)
- Distributed Representations of Words and Phrases and their Compositionality (A recent, even simpler and faster model: word2vec, Mikolov et al. 2013)
- *Glove: Global Vectors for Word Representation* (Pennington et al., 2014 and Levy and Goldberg, 2014)
- ... more later

A Simple Feed forward Neural Network

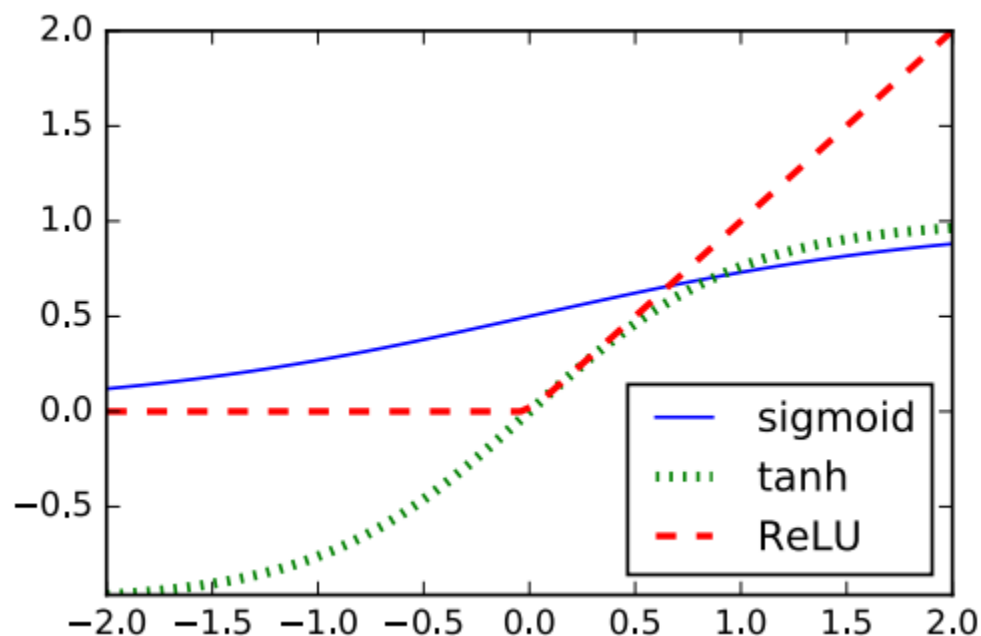
- Feed forward network



$$h = g(Vx + c)$$

$$\hat{y} = Wh + b$$

Nonlinear activation functions



$$\text{sigmoid}(x) = \frac{e^x}{1 + e^x}$$

$$\tanh(x) = 2 \times \text{sgm}(x) - 1$$

$$(x)_+ = \max(0, x)$$

a.k.a. "ReLU"

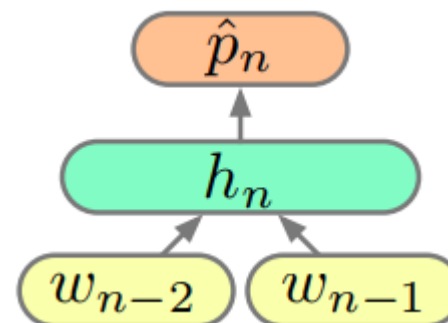
A (Very) Simple NN for Language Model



- 模拟了一个bigram语言模型
- 前一个词通过隐藏层的映射，来预测后一个词 (输出层对应于 $|V|$ 个分类器)

Trigram NN language model

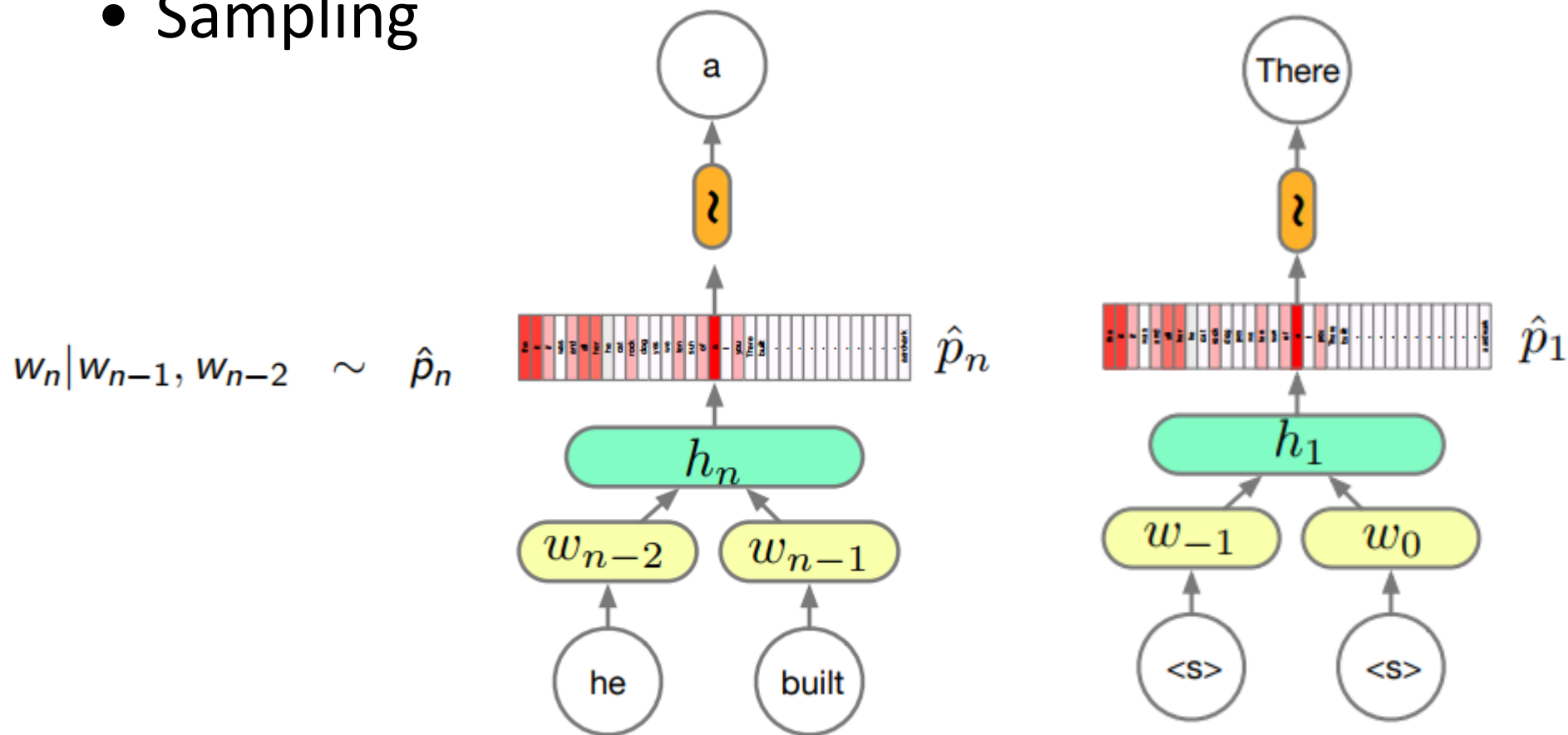
$$\begin{aligned}h_n &= g(V[w_{n-1}; w_{n-2}] + c) \\ \hat{p}_n &= \text{softmax}(Wh_n + b) \\ \text{softmax}(u)_i &= \frac{\exp u_i}{\sum_j \exp u_j}\end{aligned}$$



- w_i are one hot vectors and \hat{p}_i are distributions
- $|w_i| = |\hat{p}_i| = V$ (words in the vocabulary)
- V is usually very large $> 1e5$

Trigram NN language model

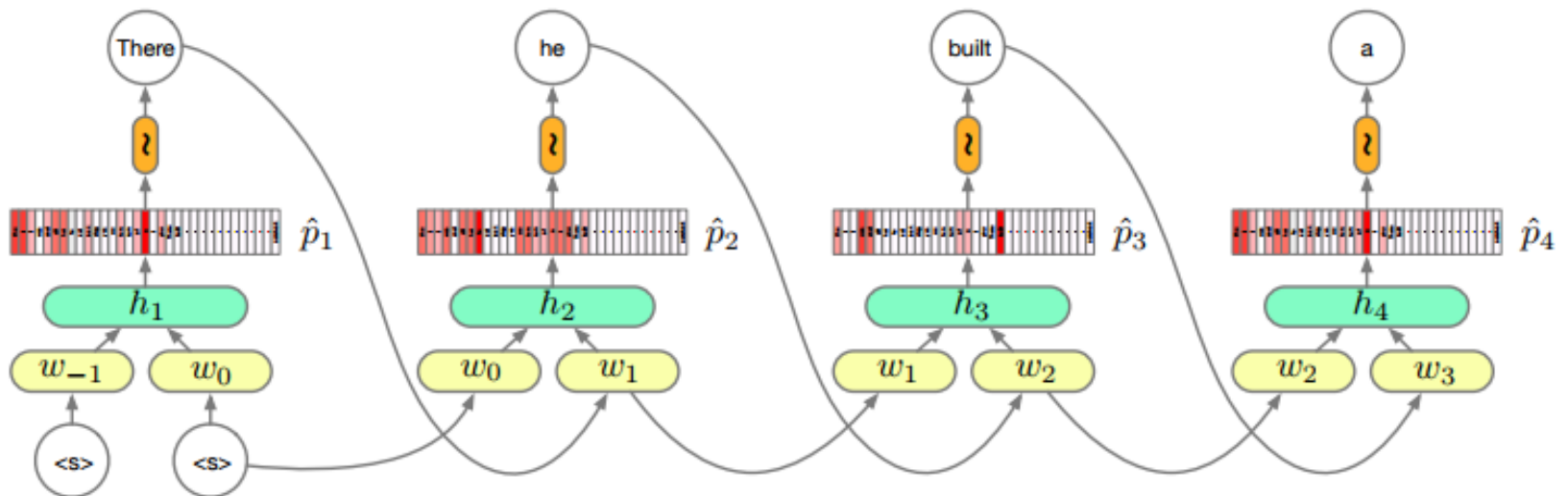
- Sampling



Trigram NN language model

- Sampling

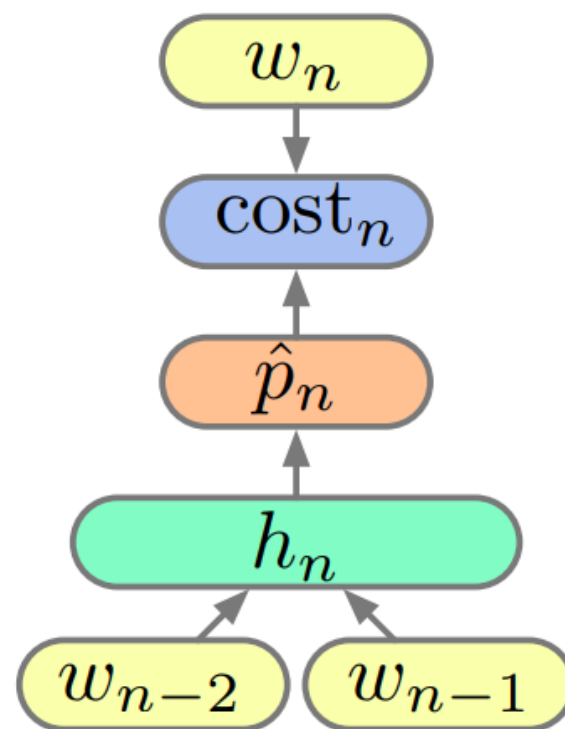
$$w_n | w_{n-1}, w_{n-2} \sim \hat{p}_n$$



Trigram NN language model

- Training: the usual training objective is the cross entropy of the data given the model

$$\mathcal{F} = -\frac{1}{N} \sum_n \text{cost}_n(w_n, \hat{p}_n)$$

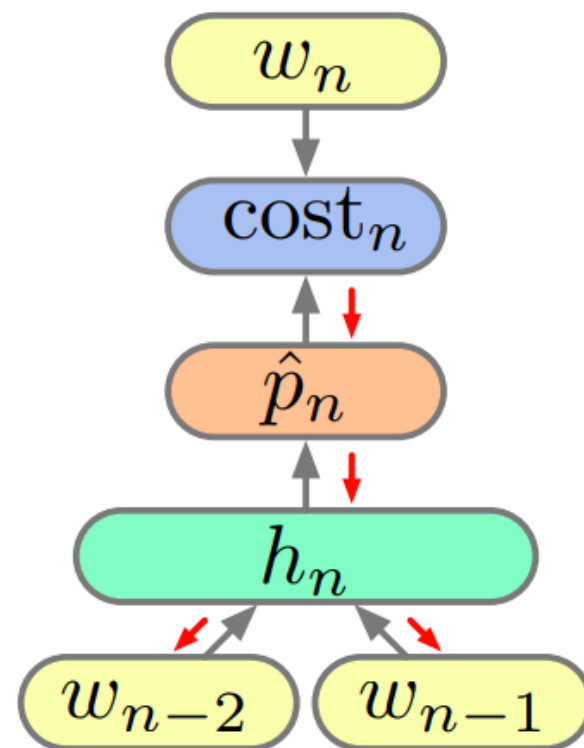


Trigram NN language model

- Training: Calculating the gradients is straightforward with back propagation

$$\frac{\partial \mathcal{F}}{\partial W} = -\frac{1}{N} \sum_n \frac{\partial \text{cost}_n}{\partial \hat{p}_n} \frac{\partial \hat{p}_n}{\partial W}$$

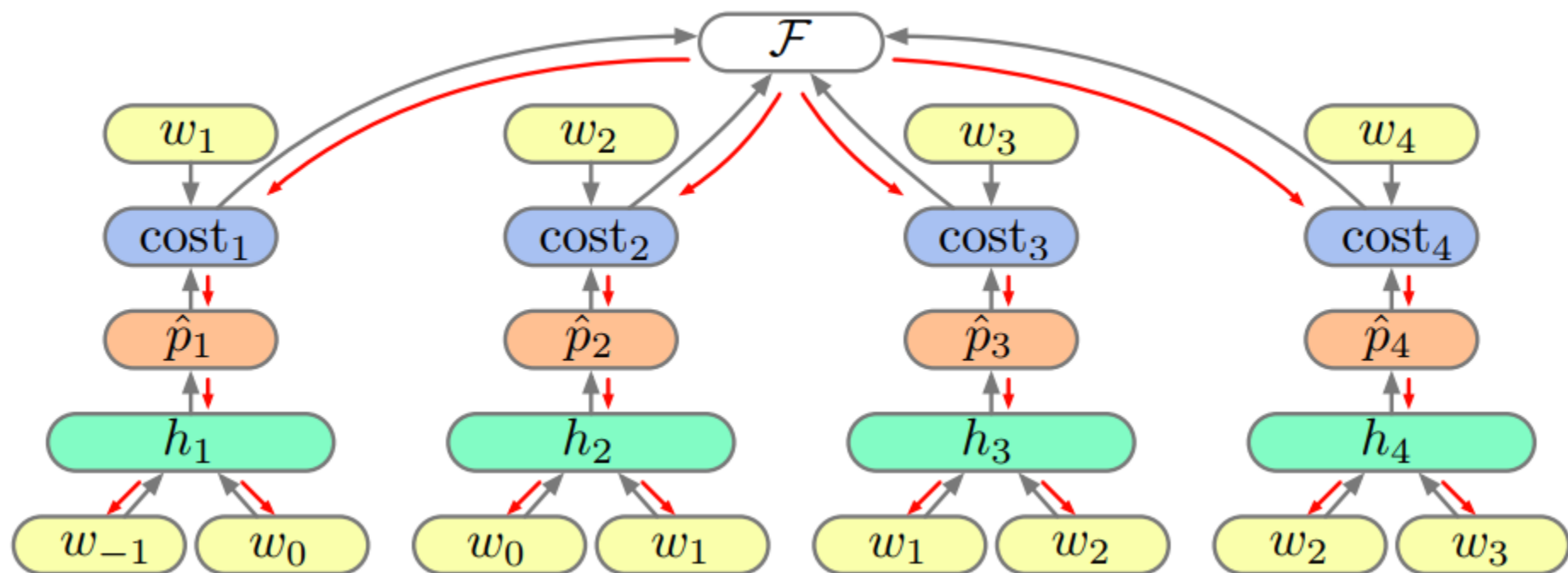
$$\frac{\partial \mathcal{F}}{\partial V} = -\frac{1}{N} \sum_n \frac{\partial \text{cost}_n}{\partial \hat{p}_n} \frac{\partial \hat{p}_n}{\partial h_n} \frac{\partial h_n}{\partial V}$$



Trigram NN language model

- Training: Calculating the gradients is straightforward with back propagation:

$$\frac{\partial \mathcal{F}}{\partial W} = -\frac{1}{4} \sum_{n=1}^4 \frac{\partial \text{cost}_n}{\partial \hat{p}_n} \frac{\partial \hat{p}_n}{\partial W} \quad \frac{\partial \mathcal{F}}{\partial V} = -\frac{1}{4} \sum_{n=1}^4 \frac{\partial \text{cost}_n}{\partial \hat{p}_n} \frac{\partial \hat{p}_n}{\partial h_n} \frac{\partial h_n}{\partial V}$$



NNLM (Bengio et. al 2003)

- Bengio et. al (2003)模型贡献
 - 得到分布式词表示
 - 得到基于分布式词表示的语言模型
 - 高阶(例如6阶), 无需平滑
 - 实验表明比基于符号的语言模型更好
 - PP值评测
- 符号约定
 - x : 变量
 - w_i : 具体词
 - v' : v 的转置
 - A_j : A 的第 j 行

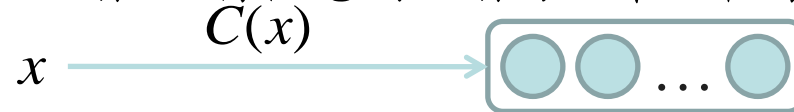
NNLM

- 模型结构:

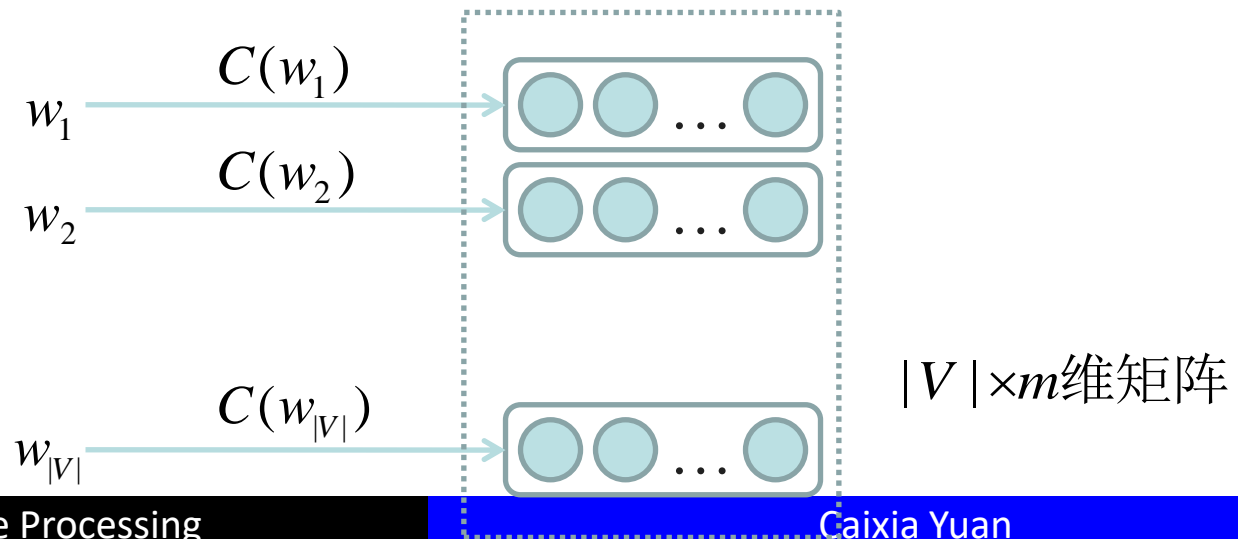
- 1. 词表映射

- 目标: 对词表 V 中的词($w_1, \dots, w_i, \dots, w_{|V|}$) 得到其 m 维向量表示
 - 实现方式

- 查表映射 C : 将任意词映射为一个 m 维向量



- 对于 V 中所有词: 将 V 中第 i 词 w_i 映射为 $C(w_i)$, 简记为 $C(i)$



• 2.神经网络模型

- 模型目标：训练一个映射 g 来建模 n 元语言模型，即

$$g(C(x_t), C(x_{t-1}), \dots, C(x_{t-n+1}); \omega) = P(C(x_t) | C(x_{t-1}), \dots, C(x_{t-n+1}))$$

- 其中 ω 为神经网络参数

- 训练的目标是使得该 n 元模型对于测试词序列 x_1, x_2, \dots, x_T (x_i 均为词表 V 中的词)具有最小PP值，即极小化：

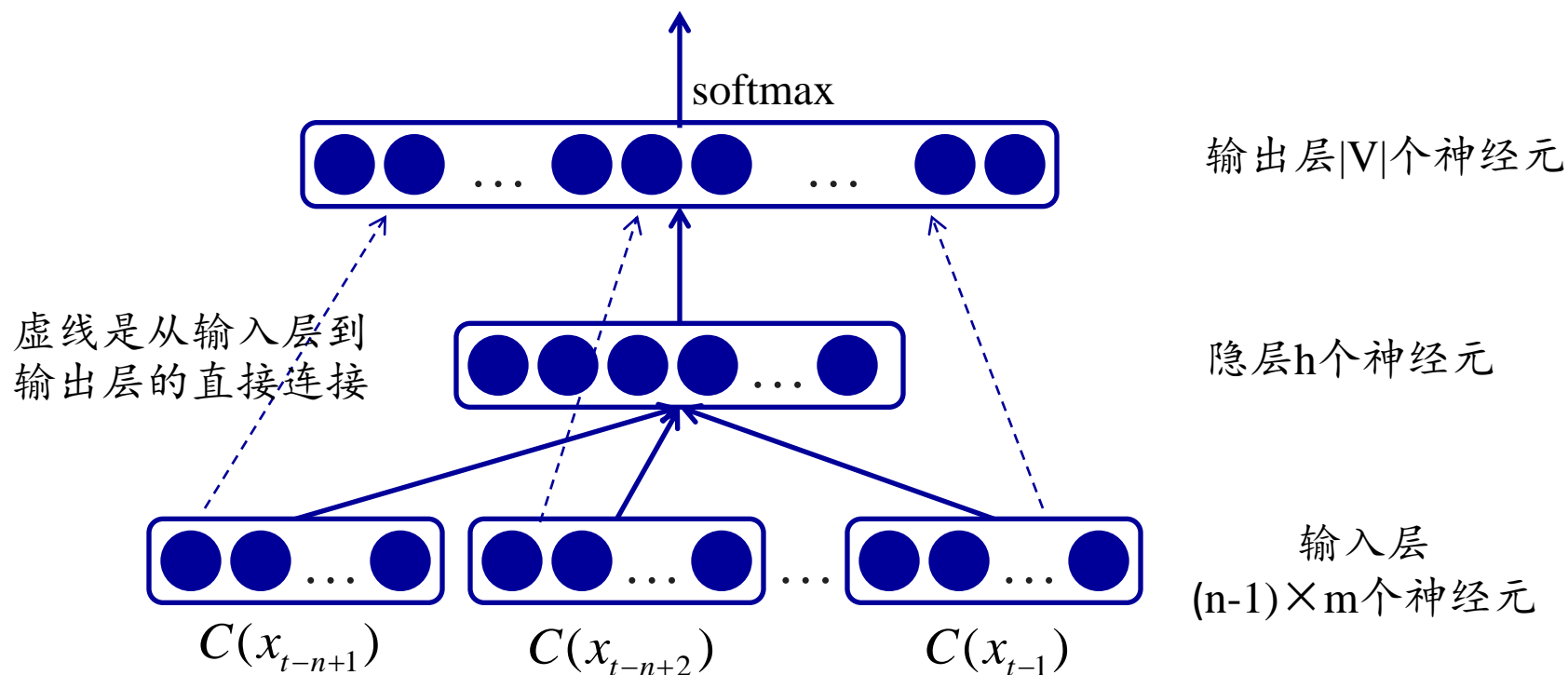
$$\begin{aligned} PP(C(x_1), \dots, C(x_T)) &= P(C(x_1), \dots, C(x_T))^{-\frac{1}{T}} \\ &= \left(\prod_{t=1}^T P(C(x_t) | C(x_{t-1}), \dots, C(x_{t-n+1})) \right)^{-\frac{1}{T}} \end{aligned}$$

- 即极大化:

$$\begin{aligned} L &= \frac{1}{T} \sum_{t=1}^T \log P(C(x_t) \mid C(x_{t-1}), \dots, C(x_{t-n+1})) \\ &= \frac{1}{T} \sum_{t=1}^T \log g(C(x_t), C(x_{t-1}), \dots, C(x_{t-n+1}); \omega) \end{aligned}$$

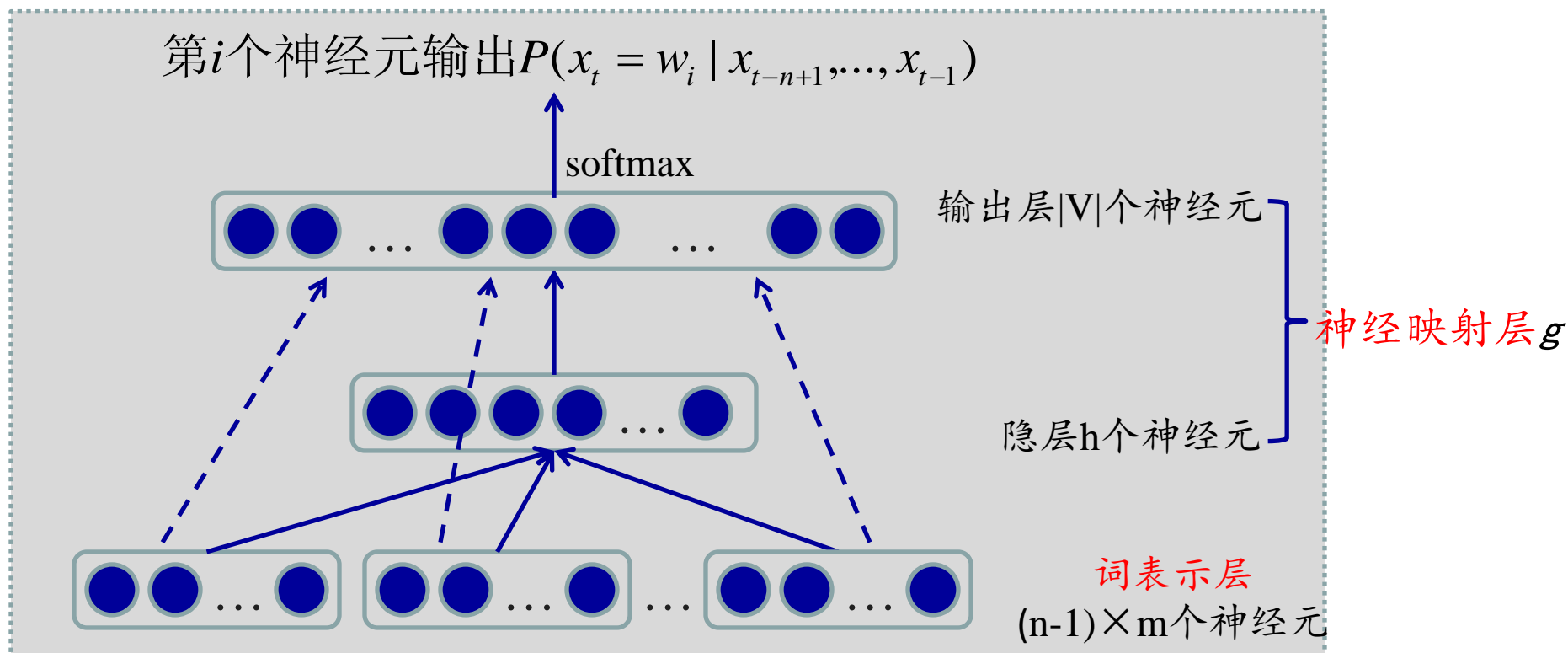
• 神经网络结构

第 i 个神经元输出为 $P(C(x_t) = C(w_i) | C(x_{t-n+1}), \dots, C(x_{t-1}))$



NNLM

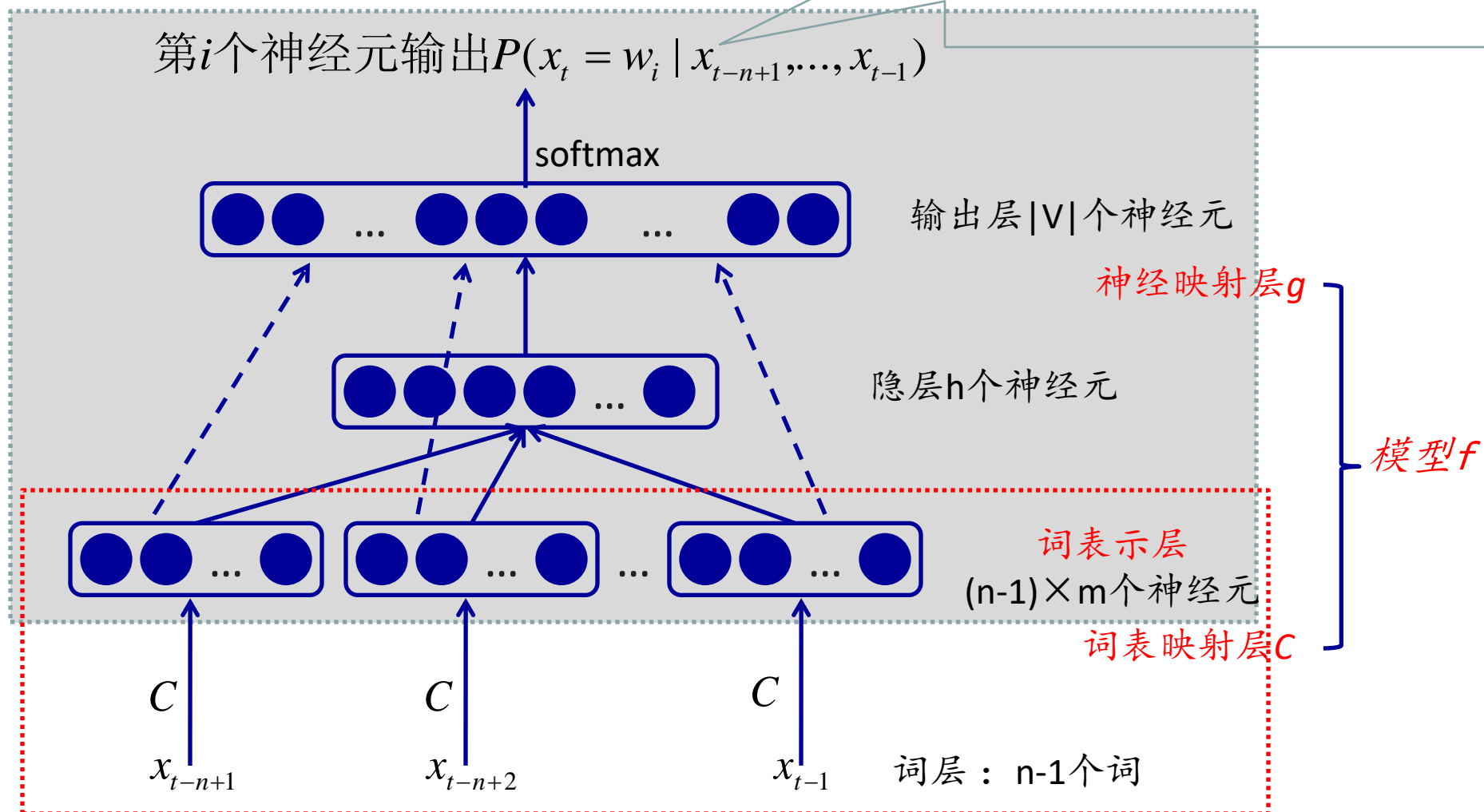
- 神经网络结构



- 合并词表映射与神经网络模型

给定C时有:

$$P(C(x_t) | C(x_{t-n+1}), \dots, C(x_{t-1})) \\ = P(x_t | x_{t-n+1}, \dots, x_{t-1})$$



- 整体模型的训练目标中把 C 也纳入，则为极大化：

$$\begin{aligned} L &= \frac{1}{T} \sum_{t=1}^T \log g(C(x_t), C(x_{t-1}), \dots, C(x_{t-n+1}); \omega) \\ &= \frac{1}{T} \sum_{t=1}^T \log f(x_t, x_{t-1}, \dots, x_{t-n+1}; C, \omega) \end{aligned}$$

- 加上正则化项，则为：

$$L = \frac{1}{T} \sum_{t=1}^T \log f(x_t, x_{t-1}, \dots, x_{t-n+1}; C, \omega) + R(C, \omega)$$

NNLM

- 模型参数

- 各层

- 词层 $n-1$ 个节点（ n -gram语法的 $n-1$ 个历史词）
 - 词表示层 $(n-1) \times m$ 个节点，每个词用 m 维向量表示
 - 隐层 h 个节点，阈值为 d ， h 维
 - 输出层 $|V|$ 个节点，阈值为 b ， $|V|$ 维

- 层间

- 词层到表示层：每一个词都对应一个向量表示，得到 $C = |V| \times m$ 矩阵
 - 表示层到隐层：权重 H ， $(n-1)m \times h$ 矩阵
 - 表示层到输出层：权重 W ， $(n-1)m \times |V|$ 矩阵
 - 隐层到输出层：权值 U ， $h \times |V|$ 矩阵

- 总参数个数

- $|V| * (1 + mn + h) + h * (1 + (n-1)m)$

NNLM

- 模型计算：对每一个输入的n元串
 - 前向计算
 - 隐层输入为： $y = b + W * C(x) + U \tanh(d + H C(x))$
 - 隐层输出为：

$$P(x_t | x_{t-1}, \dots, w_{t-n+1}) = \frac{e^{y_{x_t}}}{\sum_i e^{y_{x_i}}}$$

- 其中： $C(x)$ 是词 x 的向量表示
- 参数集： $\theta = (b, W, C, U, d, H)$
- 反向随机梯度下降

$$\theta \leftarrow \theta + \varepsilon \frac{\partial \log P(x_t | x_{t-1}, \dots, w_{t-n+1})}{\partial \theta}$$

- ε 为学习率
- 不在输入窗口中的词向量值不需要调整

Next lecture

- Language Model Evaluation
- 复习/自学：
 - 信息熵