

Natural Language Processing
#L3-4

Word Embedding & Neural Language Model

袁彩霞

yuancx@bupt.edu.cn

人工智能学院 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:

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

- 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！

基于符号的词表示

- 词的符号表示
 - 词就是一个原子符号(atomic symbol)
 - motel和hotel是不同的符号，因而是不同的词
- 这种表示被称为独热表示 (one-hot representation)
- 独热表示下，
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

基于符号的词表示

- 独热表示时一种分布式表示 (**discrete representation**)
- 在向量空间中，每个词对应着一个表示向量，其中包含一个1，和若干个0

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

- 向量的维度：

20K (speech) – 50K (PTB) – 500K (big vocabulary) – 13M (Google 1T)

如何表示词的语义

- 词的语义表示？
 - 什么是词的语义？
 - 如何表示词的语义？

如何表示词的语义

- 定义: **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)

如何表示词的语义

- 一个常用方法：通过词的类别体系区分不同词的关系，例如WordNet，表示出词的上下位关系、同义词

[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) 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

如何表示词的语义

- 作为语言资源很有价值，但是缺少词的细粒度描述，例如同义词之间的差异
 - 同义词(**synonyms**): adept, expert, good, practiced, proficient, skillful?
- 难以描述新词（保持词典同步更新十分困难）：wicked, badass, crack, ace, wizard, genius, ninja
- 主观：依赖于人类专家的语言知识
- 人工标注成本昂贵
- 词之间的相似性，难以直接度量

基于分布相似性的词的表示

- 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)

- 成为现代统计自然语言处理的一个重要思想

- e.g.,

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

- 采用上下文中的词来表示 *banking*

基于分布相似性的词的表示

- 如何使用上下文？
- 常用做法：构造词的共现矩阵 X
- 两种选择：全文 vs 上下文窗口
 - 词的上下文特定窗口中的词：同时捕捉了词的语法和语义信息
 - 词在文档中的共现：给定一个文档集合，根据词在不同文档中的出现情况，表示不同的词(和“Latent Semantic Analysis”技术关联，后续讨论)

基于窗口的共现矩阵

- 设定窗口长度 l (e.g., 1, 或更一般的5 - 10)
- 对称性(不区分上文或下文中的共现)
- 例如：给定语料
 - I like deep learning.
 - I like NLP.
 - I enjoy flying.

基于窗口的共现矩阵

- 例如：给定语料
 - 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

基于窗口的共现矩阵

- 基于共现的词向量表示存在的问题：
 - 向量维度随着词汇规模增长
 - 向量维度过大：存储代价过高
 - 稀疏问题严重 → 模型鲁棒性差

解决方法：构建词的低维向量

- 基本思想：采用一个固定长度、低维度向量(dense vector)表示词的最重要的信息
- 例如采用25 – 1000维的向量
- 如何降低维度？

Method 1: Dimensionality Reduction on X

- Singular Value Decomposition of cooccurrence matrix X

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

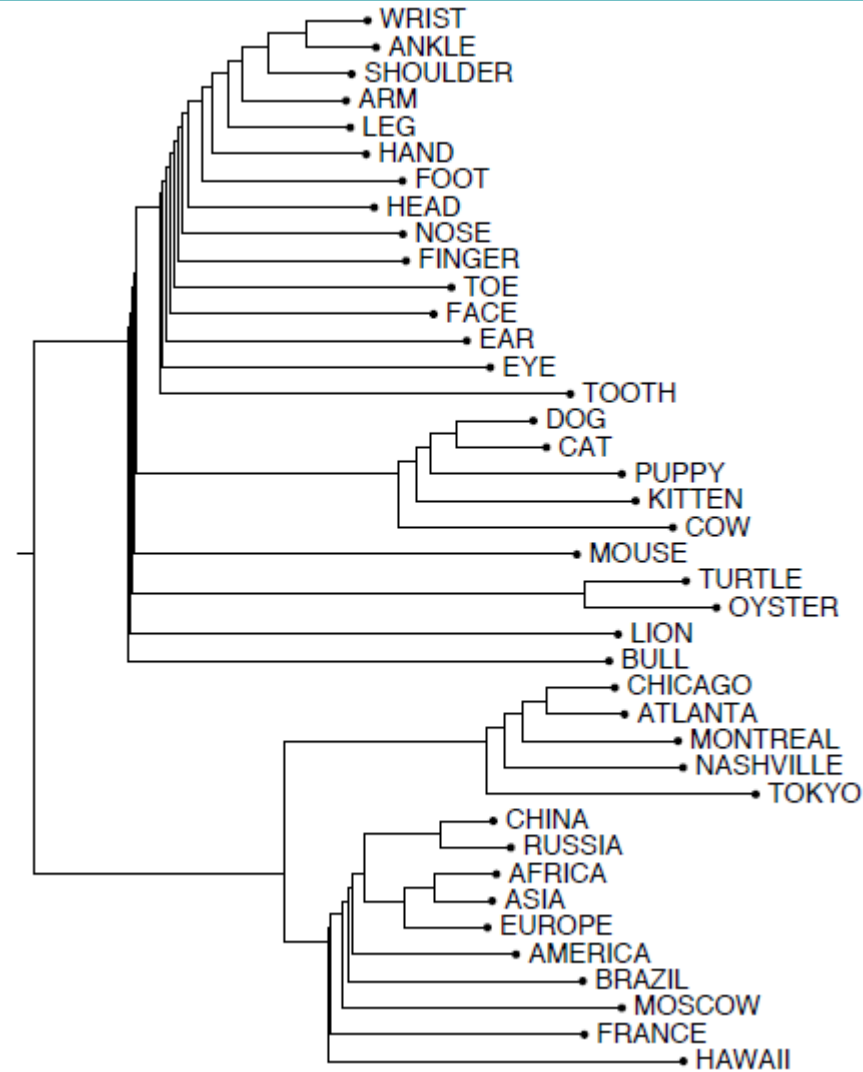
- \hat{X} is the best rank k approximation to X , in terms of least squares.

Word meaning is defined in terms of vectors

- A word is represented as a dense vector

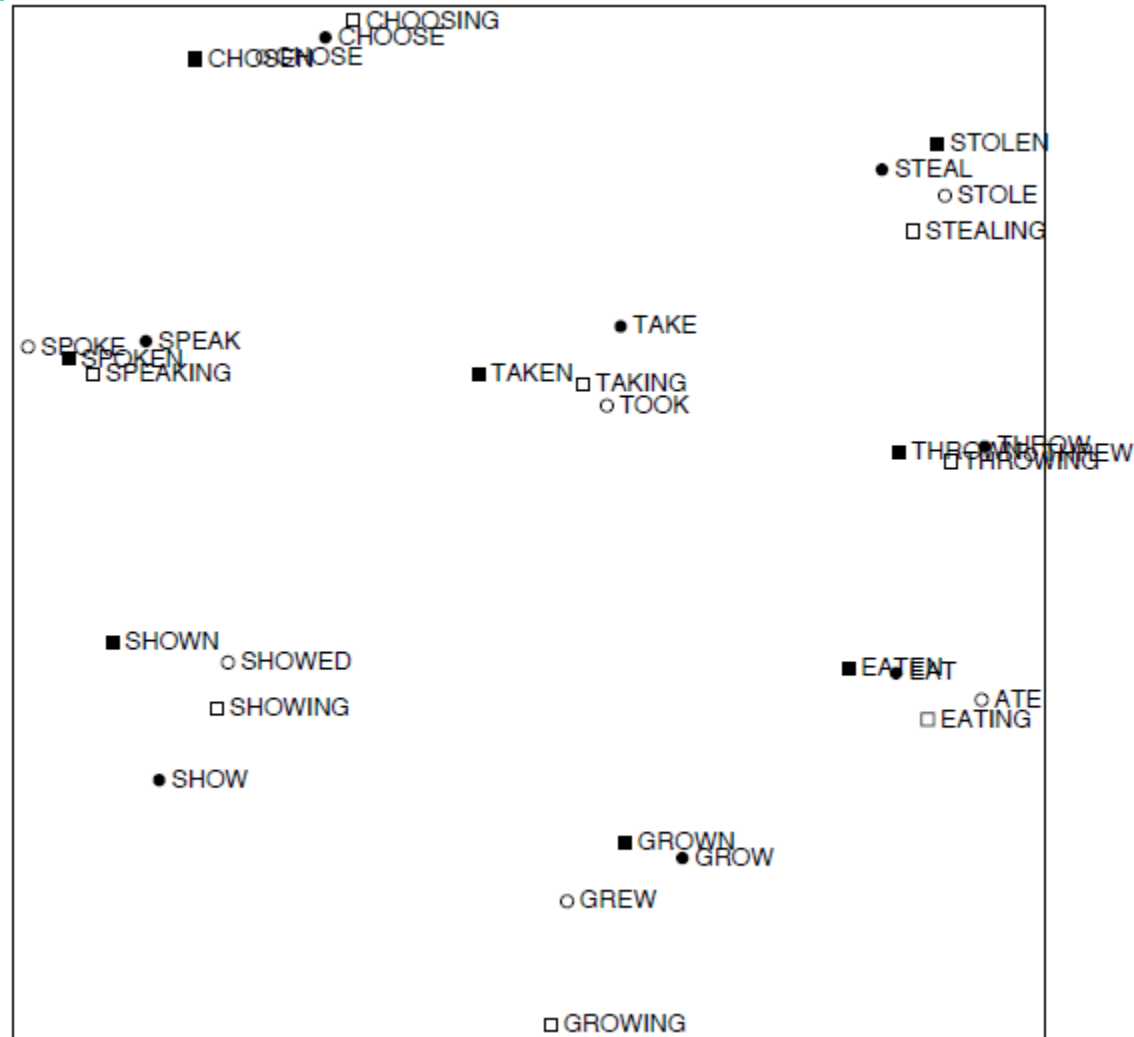
$$\textit{linguistics} = \begin{bmatrix} 0.286 \\ 0.792 \\ -0.177 \\ -0.107 \\ 0.109 \\ -0.542 \\ 0.349 \\ 0.271 \end{bmatrix}$$

Interesting semantic patterns emerge in the vectors



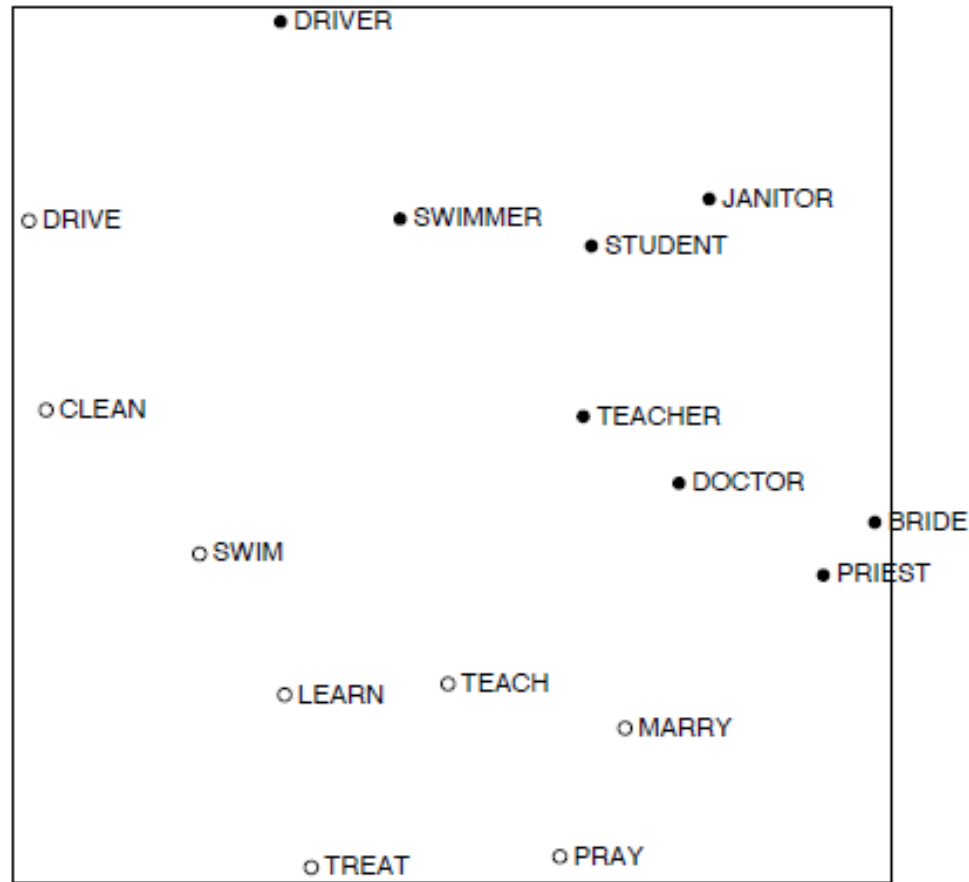
Rohde et al. An Improved Model of Semantic Similarity Based on Lexical Co-Occurrence. 2005.

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

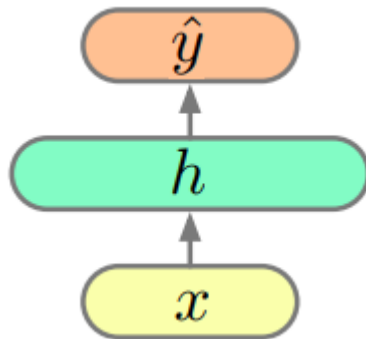
- Computational cost scales quadratically for $n \times m$ matrix: $O(mn^2)$ flops (when $n < m$)
- → Bad for millions of words or documents
- Hard to incorporate new words or documents
- Different learning regime than other deep learning models

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)
- *Deep contextualized word representation (NAACL-HLT2018)*
- *BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding* (Jacob Devlin, et al, 2018)
- ... more later

A Simple Feed forward Neural Network

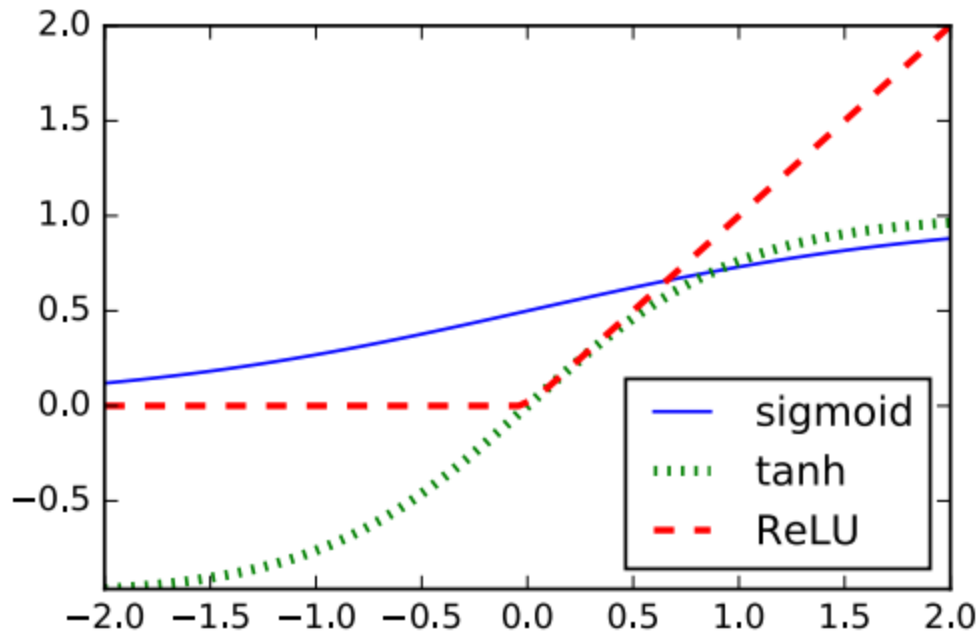
- Feed forward network



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

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

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"

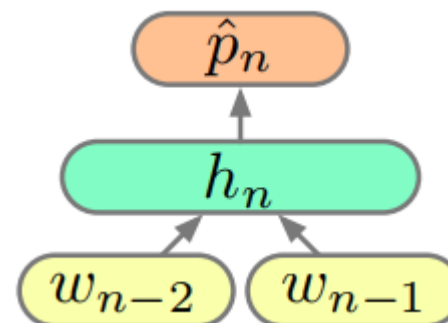
A (Very) Simple NN for Language Model



- 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

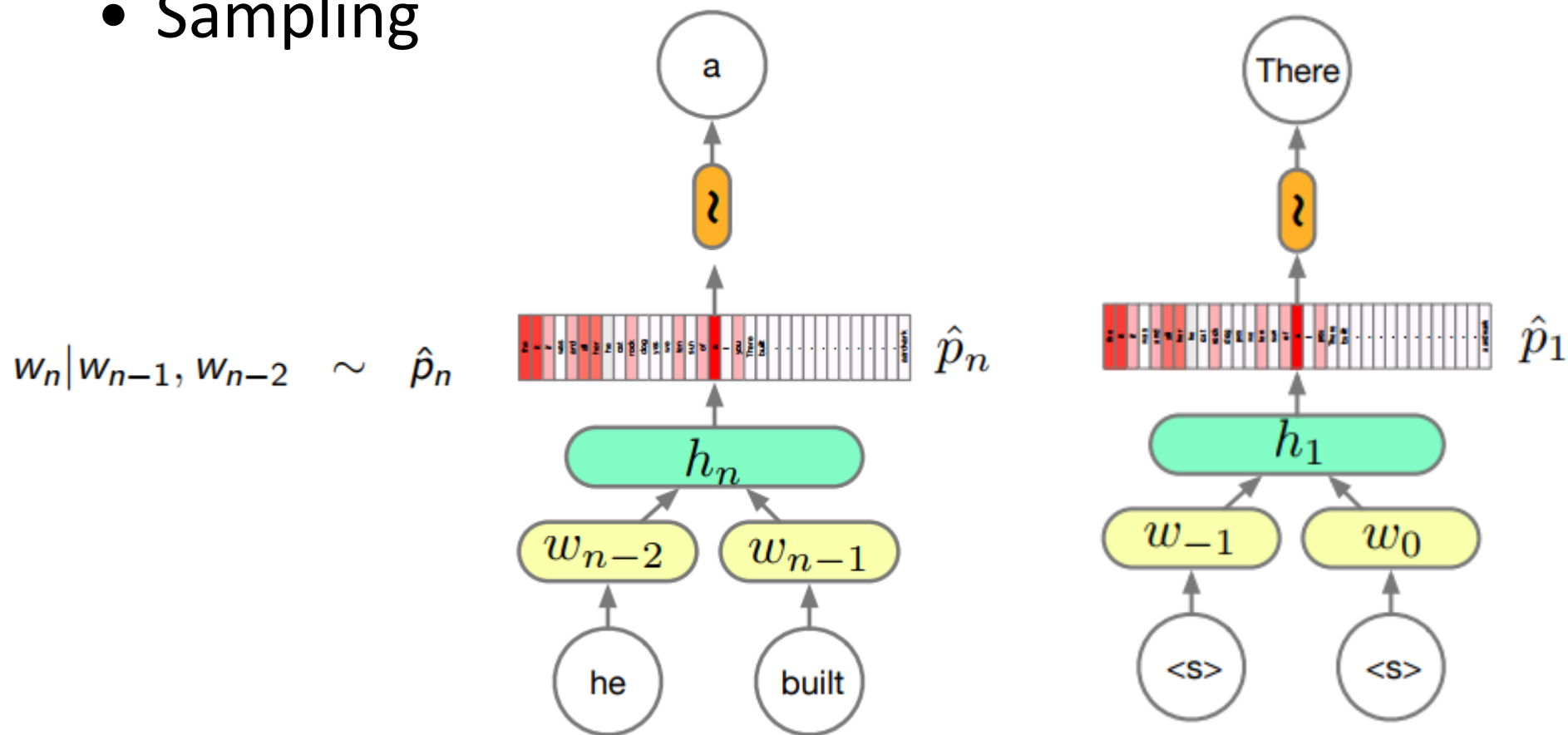
$$\begin{aligned}h_n &= g(V[w_{n-1}; w_{n-2}] + c) \\ \hat{p}_n &= \text{softmax}(Wh_n + b) \\ \text{softmax}(u)_i &= \frac{\exp u_i}{\sum_j \exp u_j}\end{aligned}$$



- 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

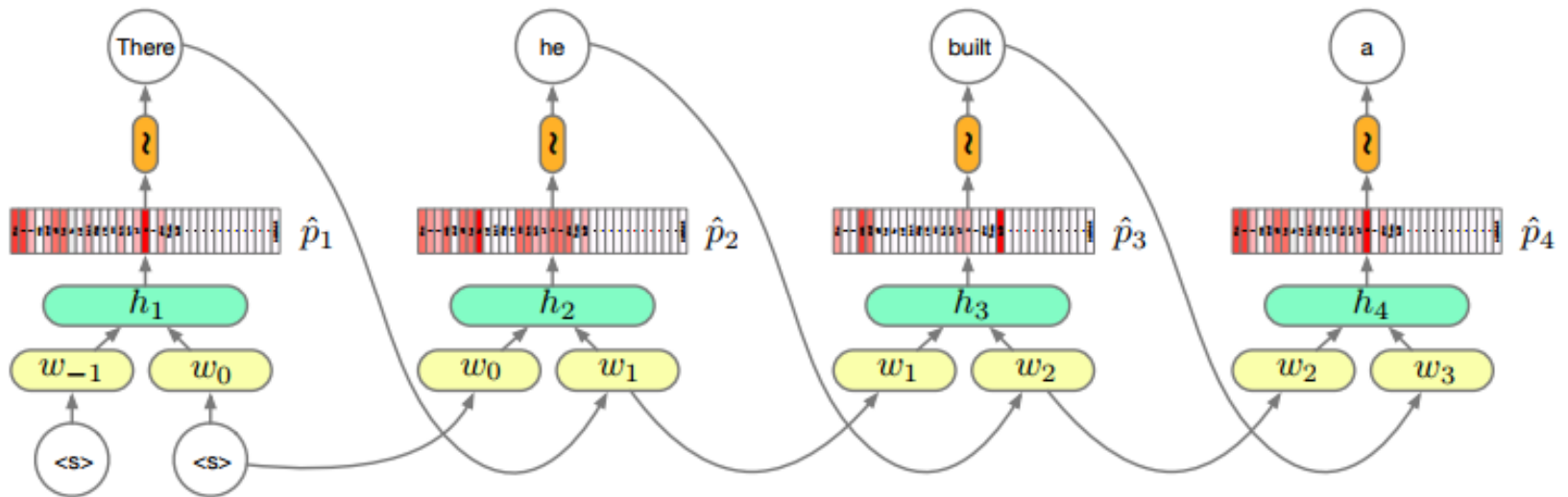
- Sampling



Trigram NN language model

- Sampling

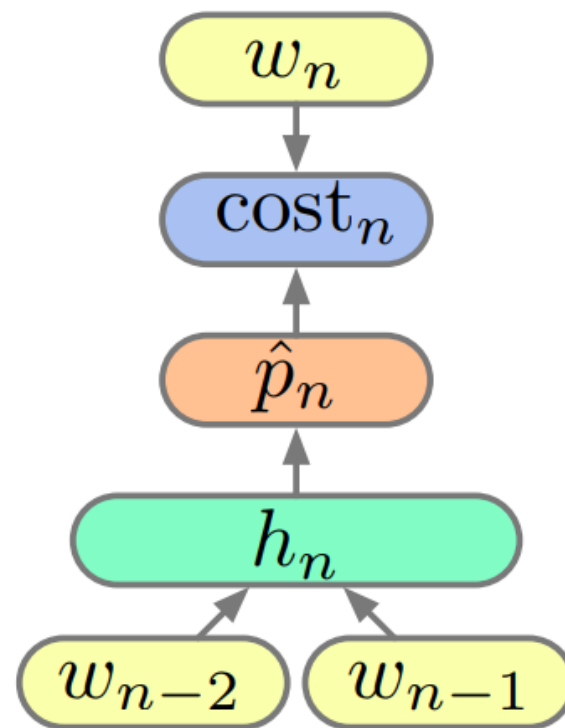
$$w_n | w_{n-1}, w_{n-2} \sim \hat{p}_n$$



Trigram NN language model

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

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

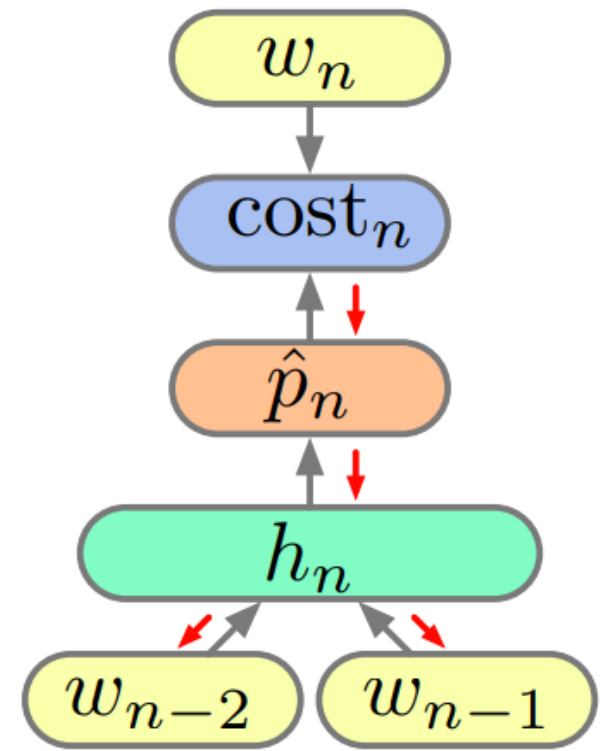


Trigram NN language model

- Training: Calculating the gradients is straightforward with back propagation

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

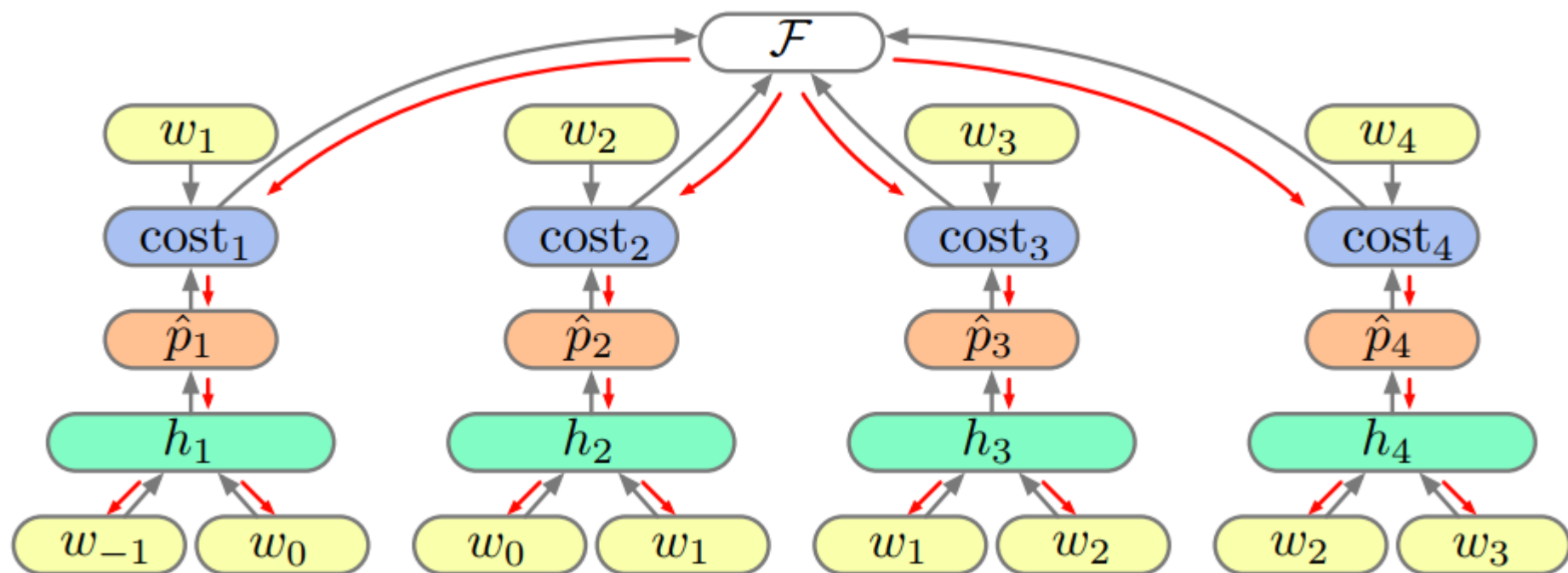
$$\frac{\partial \mathcal{F}}{\partial V} = -\frac{1}{N} \sum_n \frac{\partial \text{cost}_n}{\partial \hat{p}_n} \frac{\partial \hat{p}_n}{\partial h_n} \frac{\partial h_n}{\partial V}$$



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

- 模型结构:

- 1. 词表映射

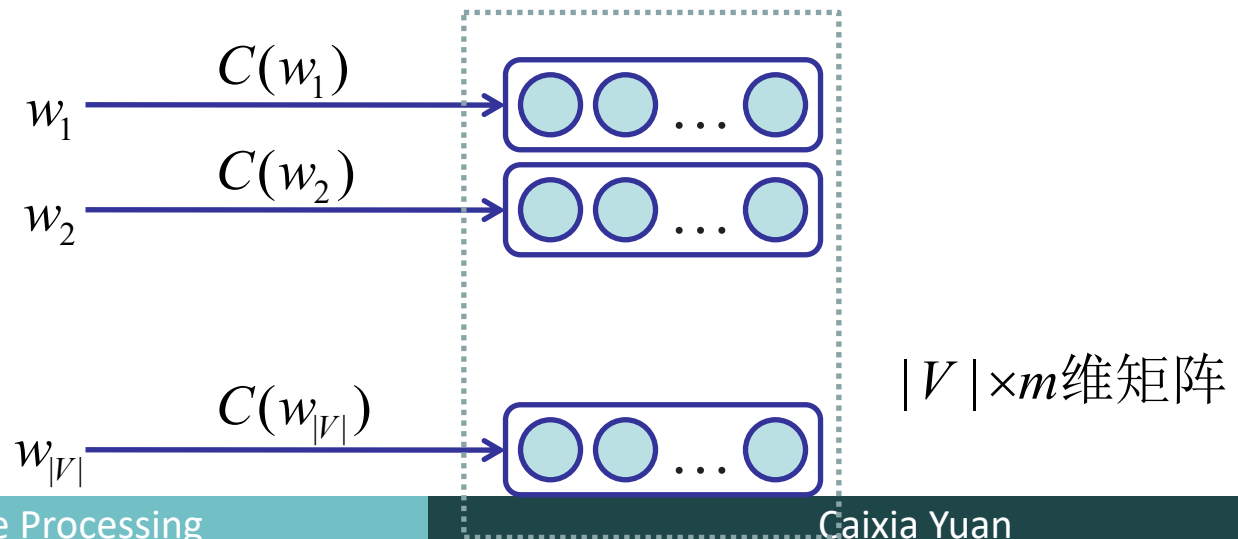
- 目标: 对词表 V 中的词($w_1, \dots, w_i, \dots, w_{|V|}$) 得到其 m 维向量表示

- 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：

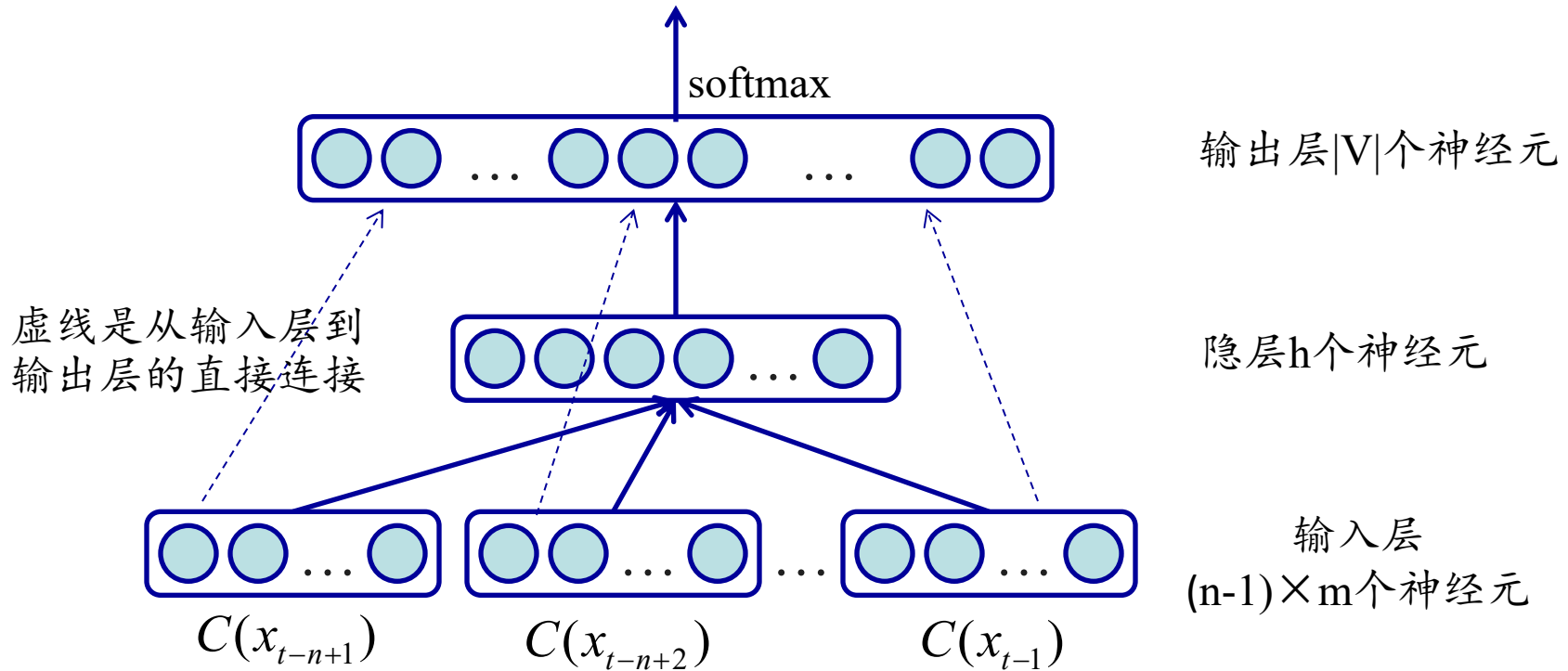
$$\begin{aligned} PP(C(x_1), \dots, C(x_T)) &= P(C(x_1), \dots, C(x_T))^{-\frac{1}{T}} \\ &= \left(\prod_{t=1}^T P(C(x_t) | C(x_{t-1}), \dots, C(x_{t-n+1})) \right)^{-\frac{1}{T}} \end{aligned}$$

- 即极大化:

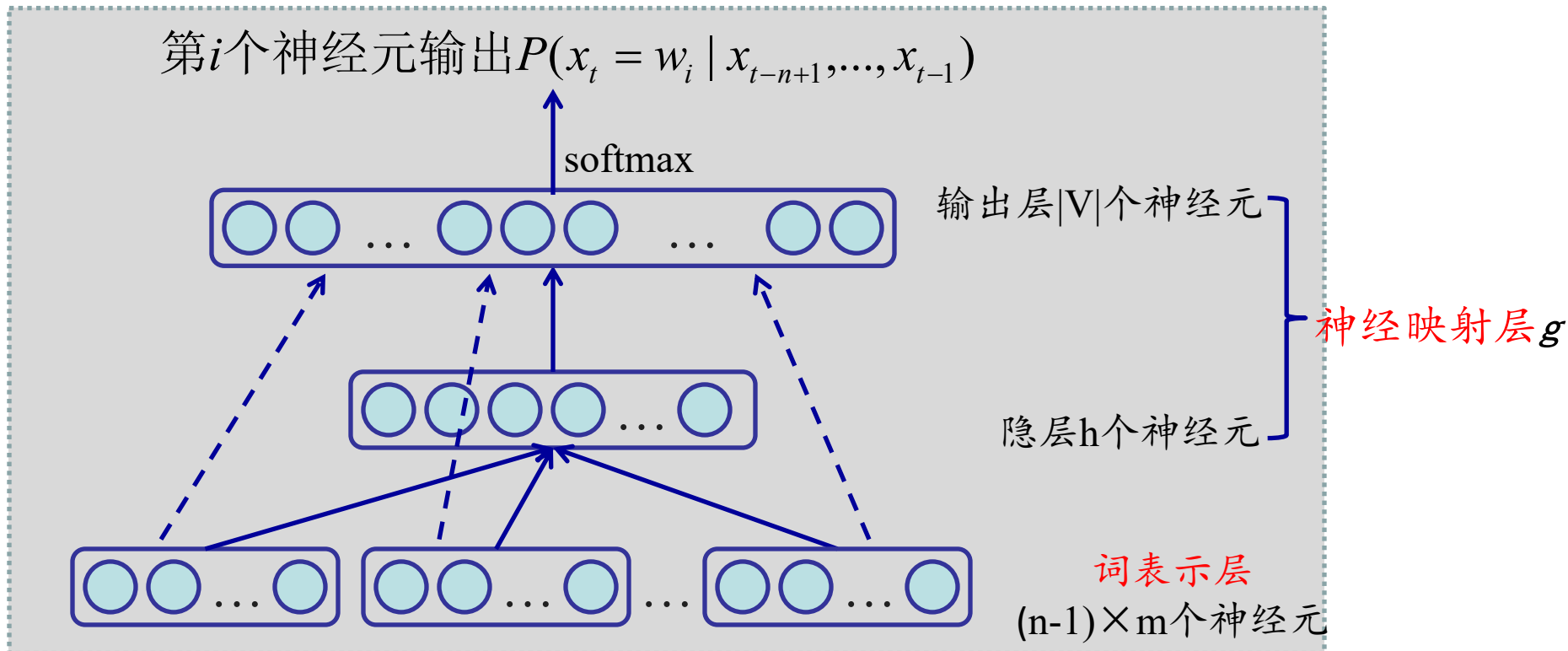
$$\begin{aligned} L &= \frac{1}{T} \sum_{t=1}^T \log P(C(x_t) | C(x_{t-1}), \dots, C(x_{t-n+1})) \\ &= \frac{1}{T} \sum_{t=1}^T \log g(C(x_t), C(x_{t-1}), \dots, C(x_{t-n+1}); \omega) \end{aligned}$$

• 神经网络结构

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



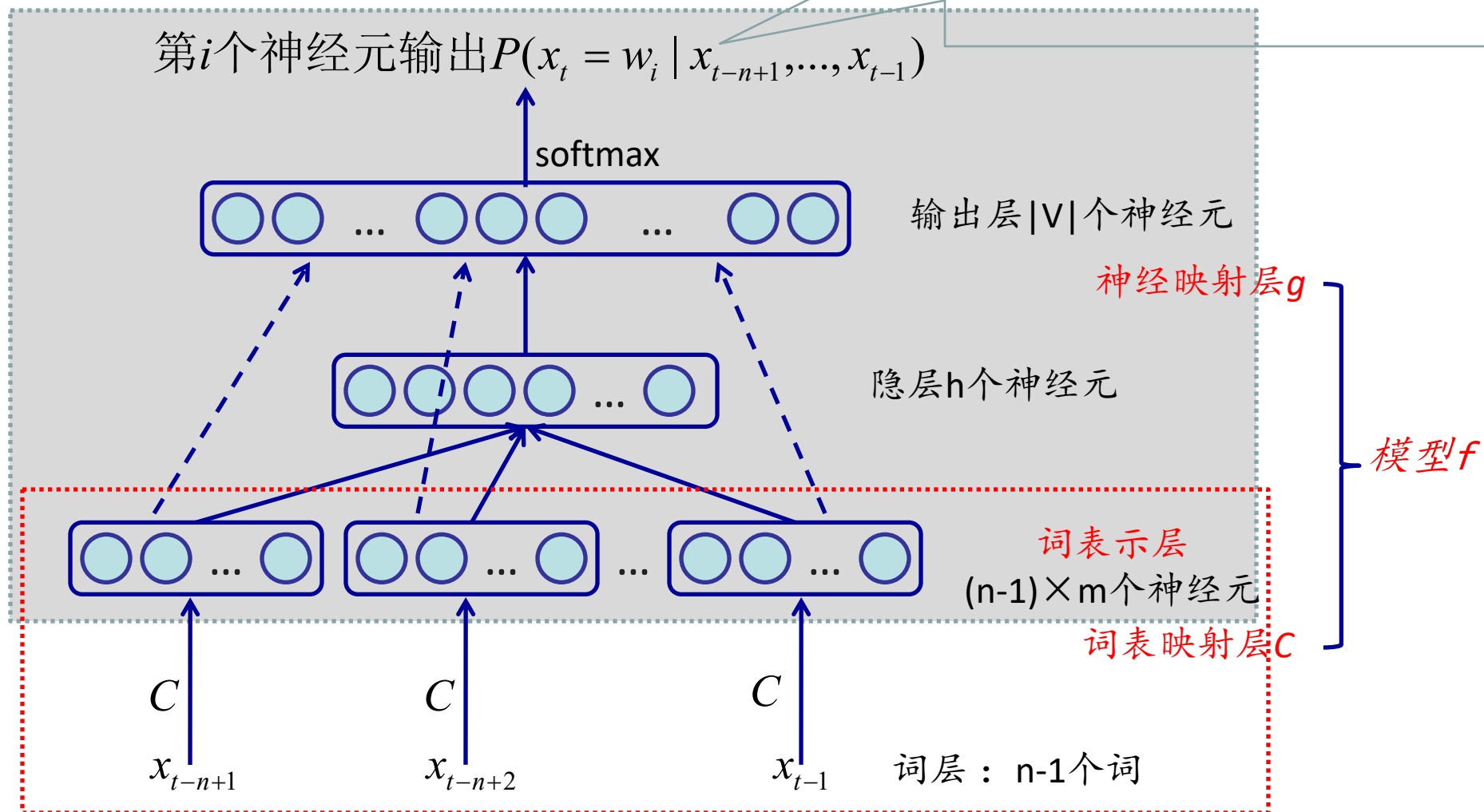
- 神经网络结构



- 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：

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

- 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

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

- 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

How the NNLM works - example

- Consider the training corpus having the following sentences:
 - “the dog saw a cat”
 - “the dog chased the cat”
 - “the cat climbed a tree”
- Vocabulary: eight words sorted alphabetically
 - “a cat chased climber dog saw the tree”
- Network architecture: eight input neurons and eight output neurons, let us assume that we decide to use three neurons in the hidden layer, and we omit bias b for simplification.
- Parameters:
 - WI and WO will be 8×3 and 3×8 matrices
 - Before training begins, these matrices are initialized to small random values as is usual in neural network training, for example,

How the NNLM works - example

$W_I =$

-0.094491	-0.443977	0.313917
-0.490796	-0.229903	0.065460
0.072921	0.172246	-0.357751
0.104514	-0.463000	0.079367
-0.226080	-0.154659	-0.038422
0.406115	-0.192794	-0.441992
0.181755	0.088268	0.277574
-0.055334	0.491792	0.263102

$W_O =$

0.023074	0.479901	0.432148	0.375480	-0.364732	-0.119840	0.266070	-0.351000
-0.368008	0.424778	-0.257104	-0.148817	0.033922	0.353874	-0.144942	0.130904
0.422434	0.364503	0.467865	-0.020302	-0.423890	-0.438777	0.268529	-0.446787

How the NNLM works - example

- Suppose we want the network to learn relationship between the words “cat” and “climbed”.
- The input vector X : $[0 \ 1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0]^T$
- The target vector: $[0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0]^T$
- The output at the hidden layer neurons can be computed as:

$$H^T = X^T W I = [-0.490796 \ -0.229903 \ 0.065460]$$

- The activation vector for output layer neurons can be computed as

$$H^T W O = [0.100934 \ -0.309331 \ -0.122361 \ -0.151399 \ 0.143463 \ -0.051262 \ -0.079686 \ 0.112928]$$

How the NNLM works - example

- Since the goal is to produce probabilities for words in the output layer, $P(word_k | word_{context})$ for $k = 1, \dots, V$
 - i.e., $\sum_k P(word_k | word_{context}) = 1$
- This can be achieved by converting activation values of output layer neurons to probabilities using the softmax function

$$y_k = \Pr(word_k | word_{context}) = \frac{\exp(activation(k))}{\sum_{n=1}^V \exp(activation(n))}$$

- Thus, the probabilities for eight words in the corpus are:
[0.143073 0.094925 0.114441 **0.111166** 0.149289 0
.122874 0.119431 0.144800]

How the NNLM works - example

- Given the target vector $[0\ 0\ 0\ 1\ 0\ 0\ 0\ 0]^T$, the error vector for the output layer is easily computed by subtracting the probability vector from the target vector.
- Once the error is known, the weights in the matrices WO and WI can be updated using backpropagation.
- In essence, this is how NNLM learns relationships between words and in the process develops vector representations for words in the corpus.

Word Analogies

- Test for linear relationships, examined by Mikolov et al. (2014)

+king [0.30 0.70]

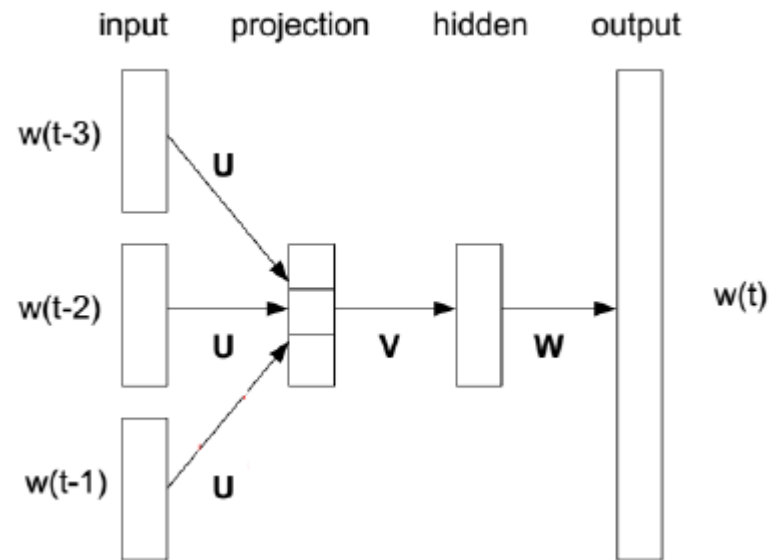
-man [0.20 0.20]

+woman [0.60 0.30]

=queen [0.70 0.80]

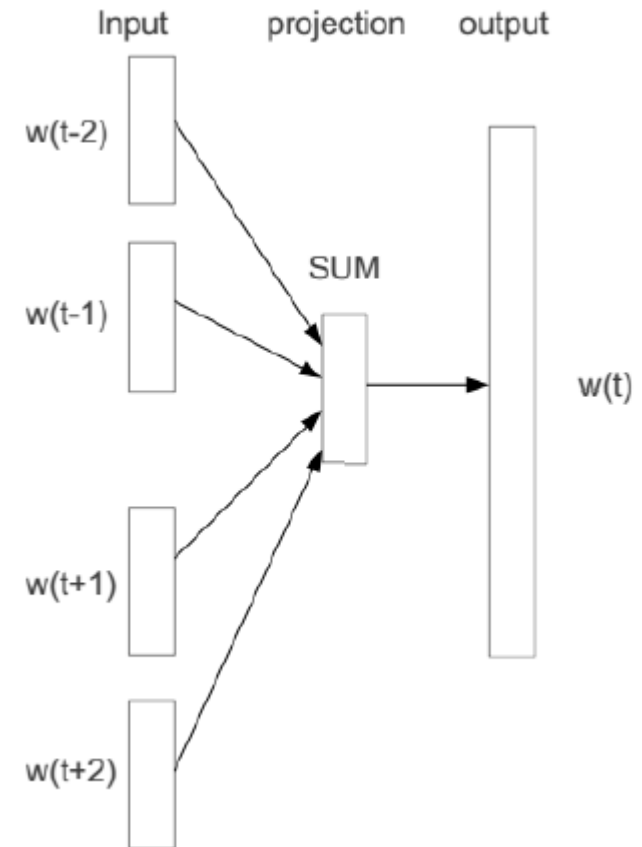
NNLM: a big picture

- Neural net based word vectors were traditionally trained as part of neural network language model (Bengio, et al, 2003)
- This models consists of input layer, projection layer, hidden layer and output layer



Other NNLMs

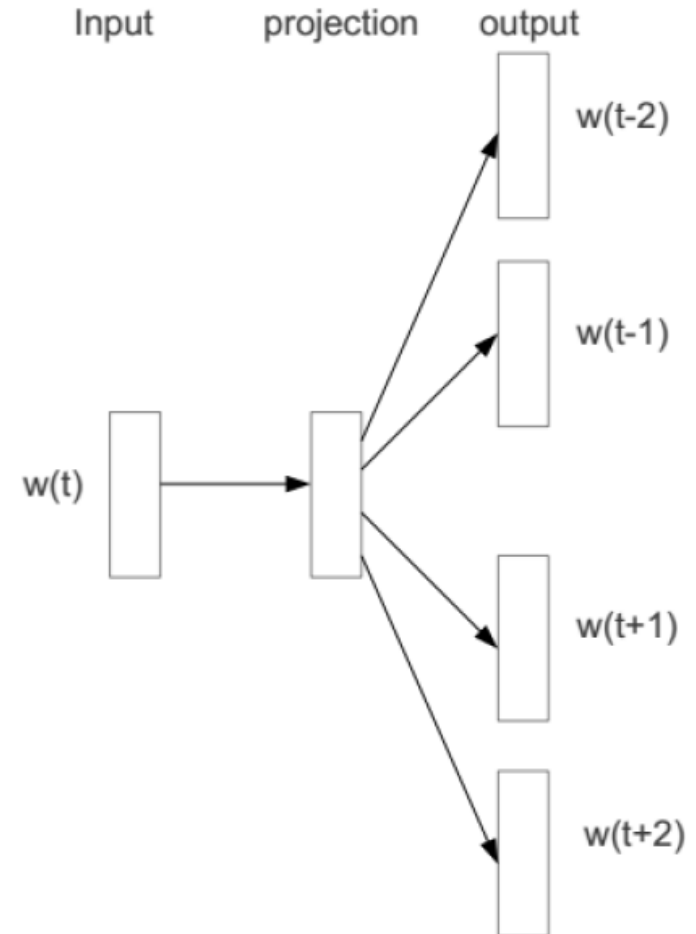
- The ‘continuous bag-of-words model’ (CBOW) adds inputs from words within short window to predict the current word
- The weights for different positions are shared
- Computationally much more efficient than normal NNLM
- The hidden layer is just linear



Mikolov et al. 2013. Distributed Representations of Words and Phrases and their Compositionality.

Other NNLMs

- We can reformulate the CBOW model by predicting surrounding words using the current word
- This architecture is called 'skip-gram NNLM'
- If both are trained for sufficient number of epochs, their performance is similar



Mikolov et al. 2013. Distributed Representations of Words and Phrases and their Compositionality.

Other NNLMs

- GloVe (Global Vectors for Word Representation)
- ELMO (Embedding from Language Models, Deep contextualized word representation)
- GPT (Generative Pre-Training)
- Transformer (Attention is all You Need)
- BERT (Bidirectional Encoder Representations from Transformers)

RNNLM toolkit

- Available at rNNLM.org
- Allows training of RNN and RNNME models
- Extensions are actively developed, for example multi-threaded version with hierarchical softmax:
`http://svn.code.sf.net/p/kaldi/code/trunk/tools/rNNLM-hs-0.1b/`

Feedforward NNLM toolkit

- Continuous Space Language Model toolkit:
<http://www-lium.univ-lemans.fr/csIm/>
- Implementation of feedforward neural network language model by Holger Schwenk

Word2vec

- Available at <https://code.google.com/p/word2vec/>
- Tool for training the word vectors using CBOW and skip-gram architectures, supports both negative sampling and hierarchical softmax
- Optimized for very large datasets (>billions of training words)
- Includes links to models pre-trained on large datasets (100B words)

Summary: NNLM

- NNLMs are currently the state-of-the-art in language modeling
- RNN outperforms FNN on language modeling tasks, both are better than n-grams, in many NLP tasks like ASR, MT
- Significant ongoing efforts to scale training to very large datasets
- The question “are neural nets better than n-grams” is incomplete: the best solution is to use both

- n-gram语言模型
 - 模型参数
 - 数据稀疏及平滑
- 神经网络语言模型
- 模型质量评价

模型质量评价

- 构建语言模型的两个阶段：
 - 模型训练：从训练数据（training data）中学习得到语言模型，即从训练样例中统计n-grams的参数，并估计其n-gram的条件概率；
 - 模型测试：在给定的测试数据（test data）中评价学习得到的语言模型
- 语言模型的建模质量如何？
- 简单说来：一个好的语言模型应该给实际使用的句子打较高的概率

实例：unigram

- 训练数据：

there is a big house
i buy a house
they buy the new house
...

- 模型：

$p(\text{there}) = 0.0714$, $p(\text{is}) = 0.0714$, $p(\text{a}) = 0.1429$
 $p(\text{big}) = 0.0714$, $p(\text{house}) = 0.2143$, $p(\text{i}) = 0.0714$
 $p(\text{buy}) = 0.1429$, $p(\text{they}) = 0.0714$, $p(\text{the}) = 0.0714$
 $p(\text{new}) = 0.0714$, ...

- 测试数据：S=they buy a big house

$$\begin{aligned} - p(S) &= 0.0714(\text{they}) \times 0.1429(\text{buy}) \times 0.0714(\text{a}) \\ &\times 0.1429(\text{big}) \times 0.2143(\text{house}) = 0.0000231 \end{aligned}$$

实例：bigram

- 训练数据：

there is a big house
i buy a house
they buy the new house
...

- 模型：

$p(\text{big} | \text{a}) = 0.5$, $p(\text{is} | \text{there}) = 1$, $p(\text{buy} | \text{they}) = 1$
 $p(\text{house} | \text{a}) = 0.5$, $p(\text{buy} | \text{i}) = 1$, $p(\text{a} | \text{buy}) = 0.5$
 $p(\text{new} | \text{the}) = 1$, $p(\text{house} | \text{big}) = 1$, $p(\text{the} | \text{buy}) = 0.5$
 $p(\text{a} | \text{is}) = 1$, $p(\text{house} | \text{new}) = 1$, $p(\text{they} | \langle s \rangle) = .333$
...

- 测试数据：S= they buy a big house

– $p(S) = 0.333 (\text{they}) \times 1 (\text{buy} | \text{they}) \times 0.5 (\text{a} | \text{buy})$
 $\times 0.5 (\text{big} | \text{a}) \times 1 (\text{house} | \text{big}) = 0.0833$

模型质量评价

- 内部评价 (Intrinsic measure) :
 - 对模型进行直接评价
- 外部评价 (Extrinsic measure) :
 - 一般用于比较两个模型A和B的质量
 - 将模型 A 和 B 应用于同一个任务 (数据)
 - 拼写检查、机器翻译、语音识别等
 - 执行该任务 (在同一个数据集中运行), 分别得到模型A和B的表现
 - 有多少个错误的词被正确检测出来?
 - 有多少个词被正确翻译?
 - 有多少个音节被正确识别?
 - 通过比较准确率来比较模型A和B的质量

困惑度

- 直接评价语言模型的质量：
 - 若测试集 W 中都是正确的句子，则 $P(W)$ 应该具有较高的概率（极大化似然）
 - 语言模型的评价指标应该是该语言模型指派给测试集的概率的函数
- Perplexity(困惑度, PP)
 - PP: $P(W)$ 的几何平均值的倒数
 - 对于测试集 W 上的n-gram模型:

$$PP(W) = P(W)^{-\frac{1}{N}} = \left[\prod_{i=1}^N P(w_i | w_{i-n}, \dots, w_{i-1}) \right]^{-\frac{1}{N}}$$

- 对于给定的测试语料 W
 - PP值越小，所用的LM越好
- 为什么？
- 解释方式1：
 - 若 W 中都是正确的句子，则 $P(W)$ 应该具有较高的概率（较小的PP值）

- 解释方式2:
 - The Shannon Game:
 - 能正确预测下一个词的能力有多强?
 - When I eat pizza, I wipe off the _____
 - Many children are allergic to _____
 - I saw a _____
 - 直观上来讲：PP指在每次预测时可能的候选词的平均个数
 - PP越大，则表明该语言的不可预测性越高
 - 语言的加权平均转移因子(weighted average branching factor)

困惑度与信息熵

- 一个词串 $W=w_1, w_2, \dots, w_N$ 的信息熵:

$$H(W) = -\sum q(w_1 \dots w_N) \log_2 q(w_1 \dots w_N)$$

- 每个词 w 的平均信息熵（熵率）:

$$H(w) = -\frac{1}{N} \sum q(w_1 \dots w_N) \log_2 q(w_1 \dots w_N)$$

- 语言 L 的熵率:

$$H(L) = -\lim_{N \rightarrow \infty} \frac{1}{N} \sum q(w_1 \dots w_N) \log_2 q(w_1 \dots w_N)$$

平稳的
遍历的

$$= -\lim_{N \rightarrow \infty} \frac{1}{N} \log_2 q(w_1 \dots w_N)$$

困惑度与信息熵

- 但 $q(.)$ 未知，在n-gram模型中，采用 p 作为 q 的估计，如果 N 足够大，则：

$$H(W) = -\frac{1}{N} \log_2 p(w_1 \dots w_N)$$

- 因此：

$$PP(W) = 2^{H(W)}$$

- H 越小，则 PP 越小，反之亦然

困惑度：实例

prediction	p_{LM}	$-\log_2 p_{\text{LM}}$
$p_{\text{LM}}(i </s><s>)$	0.109043	3.197
$p_{\text{LM}}(\text{would} <s>i)$	0.144482	2.791
$p_{\text{LM}}(\text{like} i \text{ would})$	0.489247	1.031
$p_{\text{LM}}(\text{to} \text{would like})$	0.904727	0.144
$p_{\text{LM}}(\text{commend} \text{like to})$	0.002253	8.794
$p_{\text{LM}}(\text{the} \text{to commend})$	0.471831	1.084
$p_{\text{LM}}(\text{rapporteur} \text{commend the})$	0.147923	2.763
$p_{\text{LM}}(\text{on} \text{the rapporteur})$	0.056315	4.150
$p_{\text{LM}}(\text{his} \text{rapporteur on})$	0.193806	2.367
$p_{\text{LM}}(\text{work} \text{on his})$	0.088528	3.498
$p_{\text{LM}}(. \text{his work})$	0.290257	1.785
$p_{\text{LM}}(</s> \text{work .})$	0.999990	0.000
average		2.633671

困惑度：实例

- 不同阶的n-gram模型困惑度对比：

word	unigram	bigram	trigram	4-gram
<i>i</i>	6.684	3.197	3.197	3.197
<i>would</i>	8.342	2.884	2.791	2.791
<i>like</i>	9.129	2.026	1.031	1.290
<i>to</i>	5.081	0.402	0.144	0.113
<i>commend</i>	15.487	12.335	8.794	8.633
<i>the</i>	3.885	1.402	1.084	0.880
<i>rapporteur</i>	10.840	7.319	2.763	2.350
<i>on</i>	6.765	4.140	4.150	1.862
<i>his</i>	10.678	7.316	2.367	1.978
<i>work</i>	9.993	4.816	3.498	2.394
.	4.896	3.020	1.785	1.510
</s>	4.828	0.005	0.000	0.000
average	8.051	4.072	2.634	2.251
perplexity	265.136	16.817	6.206	4.758

困惑度：实例

- 不同的平滑算法困惑度对比：

Smoothing method	bigram	trigram	4-gram
Good-Turing	96.2	62.9	59.9
Witten-Bell	97.1	63.8	60.4
Modified Kneser-Ney	95.4	61.6	58.6
Interpolated Modified Kneser-Ney	94.5	59.3	54.0

● 实验语料

	Brown语料	AP新闻
训练语料规模	1,181,041词的前800,000词	13,994,528词
发展语料规模(模型选择、权重衰减、early stopping)	随后的200,000词	963,138词
测试语料规模	其余181,041词	963,071词
语料实际含的不同词	47,578(含标点、大小写不同、分割段落与文本的标记符)	148,721词
使用的词，即 $ V $	16,383去除频率小于等于3的	17,964词(进行一些合并)
学习率	初始 $\epsilon_0=10^{-3}$ ，之后衰减，按 $\epsilon_t=\epsilon_0/(1+rt)$	
权重衰减惩罚	10^{-4}	10^{-5}
Early stopping	采用	没有
收敛	10-20epochs后	5epochs后

- 对比模型
 - Benchmark n-gram models
 - interpolated or smoothed trigram model (Jelinek and Mercer, 1980)
 - State-of-the-art n-gram models
 - back-off n-gram models with the Modified Kneser-Ney algorithm (Kneser and Ney, 1995, Chen and Goodman., 1999)
 - class-based n-gram models (Brown et al., 1992, Ney and Kneser, 1993, Niesler et al., 1998).

● Brown语料结果

模型阶数 词类数 隐单元数 词表示维数 输入到输出 是否与插值trigram混合 (权值均为0.5) PP

	n	c	h	m	direct	mix	train.	valid.	test.
MLP1	5		50	60	yes	no	182	284	268
MLP2	5		50	60	yes	yes		275	257
MLP3	5		0	60	yes	no	201	327	310
MLP4	5		0	60	yes	yes		286	272
MLP5	5		50	30	yes	no	209	296	279
MLP6	5		50	30	yes	yes		273	259
MLP7	3		50	30	yes	no	210	309	293
MLP8	3		50	30	yes	yes		284	270
MLP9	5		100	30	no	no	175	280	276
MLP10	5		100	30	no	yes		265	252
Del. Int.	3						31	352	336
Kneser-Ney back-off	3							334	323
Kneser-Ney back-off	4							332	321
Kneser-Ney back-off	5							332	321
class-based back-off	3	150						348	334
class-based back-off	3	200						354	340
class-based back-off	3	500						326	312
class-based back-off	3	1000						335	319
class-based back-off	3	2000						343	326
class-based back-off	4	500						327	312
class-based back-off	5	500						327	312



- 初步结论：
 - 更多上下文时(高阶语言模型)神经模型性能改善，而原有LM没有太多受益
 - NN模型的隐单元(有无以及数量变化)是有影响的
 - NN模型与原有LM的混合是有帮助的
 - 从输入到输出的直接连接是否有用(图中看不出，但是作者有如下讨论):
 - 小语料时提供更好的泛化能力(很有限)，大语料时直接连接提供更快的收敛速度(2倍)

- AP语料结果

	n	h	m	direct	mix	train.	valid.	test.
MLP10	6	60	100	yes	yes		104	<u>109</u>
Del. Int.	3						126	132
Back-off KN	3						121	127
Back-off KN	4						113	119
Back-off KN	5						112	117

- 由于语料规模大，所以只执行了5epochs迭代，结论和前面相似。

其它的语言模型

- Syntactic language models: using parse trees
- Class-based N-gram Model
- Topic-based N-gram Model
- Skip n-gram models: back-off to $p(w_n | w_{n-2})$

- 补充知识：
 - 似然函数及最大似然估计
 - 几种不同的实验设置

Likelihood

- In statistics, a **likelihood function** (often simply the **likelihood**) is a function of the parameters of a statistical model given data.
- The *likelihood* of a set of parameter values, θ , given outcomes x , is equal to the *probability* of those observed outcomes given those parameter values, that is

$$L(\theta|x)=P(x | \theta)$$

- The likelihood function is defined differently for discrete and continuous probability distributions.

Likelihood

- **Discrete probability distribution**
- Let X be a random variable with a discrete probability distribution p depending on a parameter θ . Then the function

$$L(\theta|x)=P(x=X | \theta)$$

- considered as a function of θ , is called the likelihood function (of θ , given the outcome x of the random variable X).

Maximum likelihood estimation

- 最大似然估计提供了一种给定观察数据来评估模型参数的方法
 - 即：“模型已定，参数未知”
- 最大似然估计的假设：所有采样独立同分布
- 例： x_1, \dots, x_n 为独立同分布的采样， θ 为模型参数， f 为我们所使用的模型
- 则参数为 θ 的模型 f 可表示为：

$$f(x_1, x_2, \dots, x_n | \theta) = f(x_1 | \theta) \times f(x_2 | \theta) \dots f(x_n | \theta)$$

- 似然函数定义为：

$$L(\theta | x_1, \dots, x_n) = f(x_1, \dots, x_n | \theta) = \prod_{i=1}^n f(x_i | \theta)$$

Maximum likelihood estimation

- 实际应用中常对两边取对数，得到：

$$\ln L(\theta | x_1, \dots, x_n) = \sum_{i=1}^n \ln f(x_i | \theta) \quad \hat{\ell} = \frac{1}{n} \ln L$$

- 其中 $\ln L(\theta | x_1, \dots, x_n)$ 称为对数似然，而 $\hat{\ell}$ 称为平均对数似然
- 最大对数平均似然，得到：

$$\hat{\theta}_{mle} = \arg \max_{\theta \in \Theta} \hat{\ell}(\theta | x_1, \dots, x_n)$$

为什么是MLE?

- 将一枚硬币抛n次，观测到的结果可以表示为概率p（p未知）的函数：

$$\begin{aligned}P(DATA | p) &= P(X_1, \dots, X_n | p) = \binom{n}{k} \prod_{i=1}^n P(X_i | p) \\&= \binom{n}{k} p^k (1-p)^{n-k} \propto \log p^k (1-p)^{n-k} \\&= k \log p + (n-k) \log(1-p) \equiv f(p)\end{aligned}$$

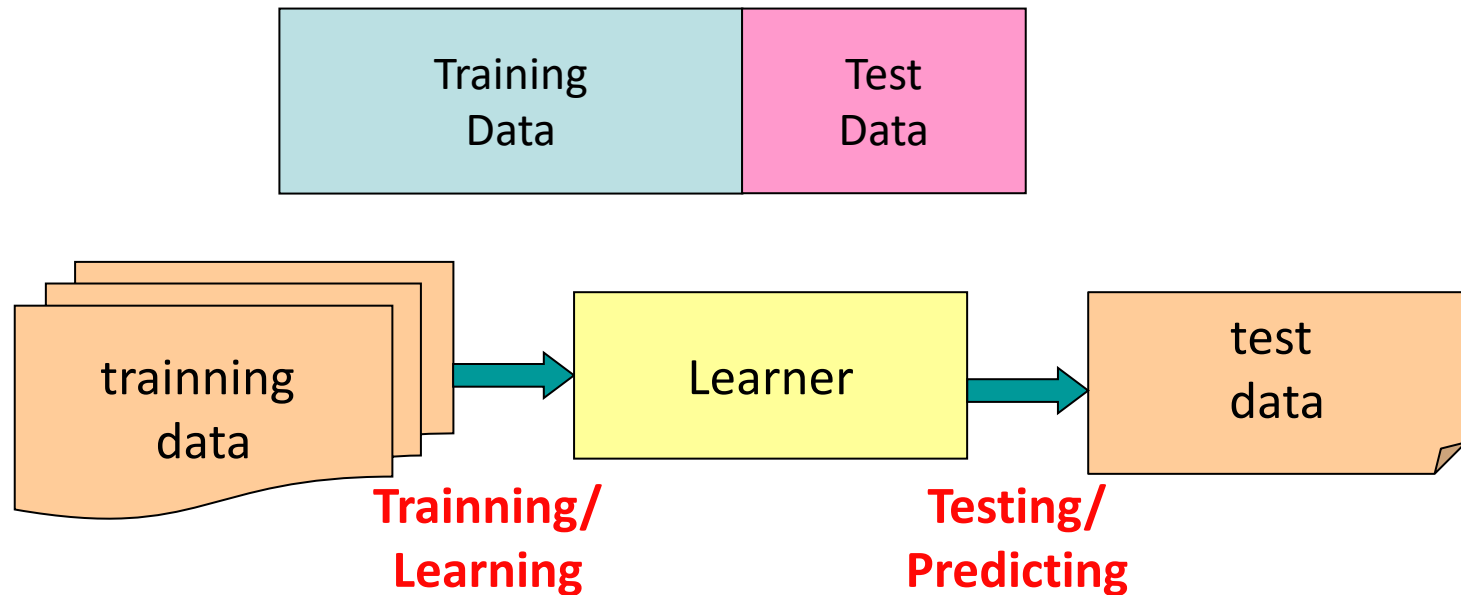
- 令 f 相对于 p 的一阶导数为零，则：

$$0 = \frac{d}{dp} f(p) = \frac{k}{p} - \frac{n-k}{1-p}, \quad (1-p)k = (n-k)p, \quad p = \frac{k}{n}$$

几种不同的实验设置

- **Setting 1: training + testing**

- 将数据集分为两部分：训练集和测试集
- 在训练集上进行模型学习和训练
- 在测试集上进行模型评估

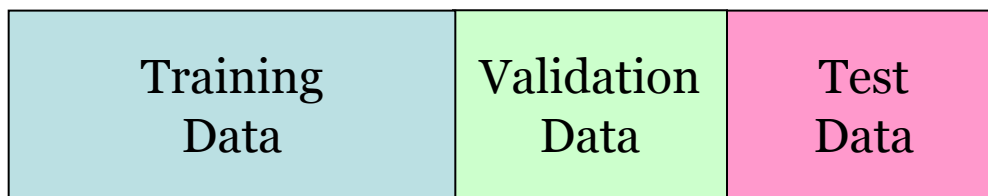


Test errors give an honest assessment of the error for future cases!

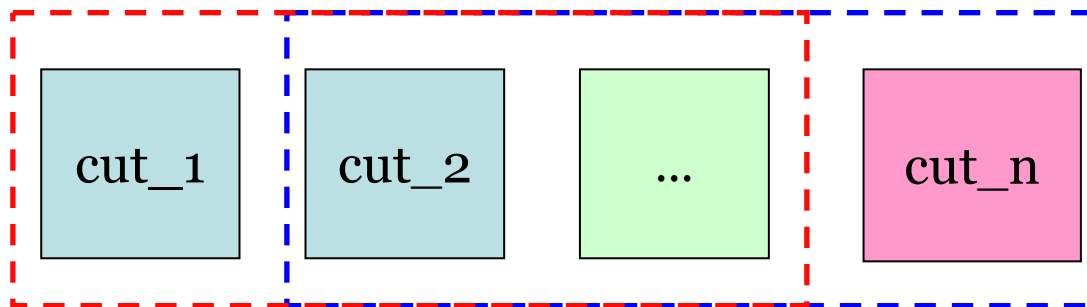
几种不同的实验设置

- **Setting 2: training + validation + testing**

- 将数据集分为三部分：训练集、验证集(发展集)、测试集
- 在训练集上进行模型训练
- 在验证集(发展集)上调整模型(超)参数
- 在测试集上进行模型评估



- **Setting 3: k倍交叉验证(k-fold Cross-Validation)**



Next lecture

- Part-of-speech
- Hidden Markov Model