



神经概率语言模型

Neural Probabilistic Language Model

王小捷
智能科学与技术中心
北京邮电大学

大纲

■ 引言

■ 词表示

■ 神经概率语言模型 (NPLM)

■ 总结





引言

■ N-gram语言模型中: $P(w_1^T) = \prod_{t=1}^T P(w_t | w_{t-n+1}^{t-1})$

■ 两个问题

■ 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(<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?



■ LM的一个缺陷：建模时没有考虑词的相似性

■ XX **cat** **is** XX 的训练样本对计算 $p(\text{is}/\text{dog})$ 无贡献

■但是，由于cat和dog的相似性， $p(\text{is}/\text{dog})$ 很可能与 $P(\text{is}/\text{cat})$ 相似

■改进：基于类的LM： $p(v/w) \approx p(v/c)p(c/w)$

■此时：如果 w_1 和 w_2 同属于类 c ，则 $p(v/w_1) \approx p(v/w_2)$

■更一般地，希望有如下性质：

■如果 x, y 是两个很相似的词，则 $f(x)$ 接近 $f(y)$ 。

■如果： $f(x)$ 是连续函数， x 是连续变量，则这是一个常见的性质

■ $f(x)$ 是连续函数： $f(x)$ 是概率，可以构建成连续函数

■ x 是连续变量： x 是词，词的连续表示？

大纲

■ 引言

■ 词表示

■ 神经概率语言模型 (NPLM)

■ 总结





■词表示

■符号表示

■cat、dog、table...

■数值表示：一个词用一个多维向量表示

■One-hot：向量维数=词表大小

■ $(1,0,0,\dots), (0,1,0,\dots)\dots$

■依靠词自身来表示词，和符号表示没有本质差异



■词表示

■符号表示

■cat、dog、table...

■数值表示：一个词用一个多维向量表示

■One-hot：向量维数=词表大小

■ $(1,0,0,\dots), (0,1,0,\dots)\dots$

■Distributional：向量维数=词表大小

■布尔式、频率式 →



■ Distributional 词表示

■ 语料C, 词表 $W=\{w_1, w_2, \dots, w_N\}$, 窗口k

■ $w \in W$ 的词表示: 用C中所有w的(窗口为k的)上下文中的词来构造w的表示:

■ 布尔式: $(I_{w1}, I_{w2}, \dots, I_{wN})$

■ 频率式: $(f_{w1}, f_{w2}, \dots, f_{wN})$

■ 例: $C=\{I \text{ am here}, I \text{ am fine.}\}$, $W=\{I, \text{am}, \text{here}, \text{fine}\}$, $k=1$

■ I

■ 布尔式: (0,1,0,0)

■ 频率式: (0,2,0,0)

■ am

■ 布尔式: (1,0,1,1)

■ 频率式: (2,0,1,1)

■

One-hot

I (1,0,0,0)

am (0,1,0,0)

here (0,0,1,0)

fine (0,0,0,1)



■ 无论是one-hot 还是distributional

■ 离散表示

■ 维度高，词表维度，每一个维度表示一个词，

- one-hot：该词本身

- distributional：该词上下文中的词

■ 低维连续表示

■ 基于one-hot 或distributional 的降维表示

- SVD分解等，后续LDA模型时再提

■ 词向量(distributed)



■词表示

■符号表示

- cat、dog、table...

■分布表示：一个词用一个多维向量表示

- One-hot：向量维数=词表大小

- (1,0,0,...), (0,1,0,...)...)

- Distributional：向量维数=词表大小

- 布尔式、频率式

- Distributed：向量维数=指定大小

- 连续值：从PNLM等模型训练获得



■ Distributed 词表示

■ 词表 $W = \{w_1, w_2, \dots, w_N\}$

■ m 维词表示

■ $c(w_1) = (x_{11}, x_{12}, \dots, x_{1m})$

■ $c(w_2) = (x_{21}, x_{22}, \dots, x_{2m})$

■ ...

■ $c(w_N) = (x_{N1}, x_{N2}, \dots, x_{Nm})$

■ 表示容量没问题：即使 x_{ij} 为 $\{0, 1\}$ ， m 也可以远远小于 N

■ 实际上 x_{ij} 为实数



■ Distributed 词表示

■ 设 $C(w) = (x_1, x_2, \dots, x_m)$ 是词 w 的 Distributed 表示

■ 则每个词是 m 维空间中的一个点

■ 希望：相似的词具有距离相近的词向量

■ 例如：牛、马、网：对于机器没有距离差别

■ 如果：

■ $c(\text{牛}) = (0.5, 0.4, 0.20)$

■ $c(\text{马}) = (0.5, 0.4, 0.09)$

■ $c(\text{网}) = (0.2, 0.7, 0.18)$

■ 则很容易计算出不同词之间的相近程度不同



■词向量用于语言模型的好处：

- 将语言模型构建为关于连续词表示的连续函数，
则可以在语言模型中充分利用词的相似性获得更好的语言模型，即相似的词(元组)得到相似的概率：
- 例如：以bigram 为例：
- 1、将 $p(w|x)$ 建模为连续函数
- 2、如果 x 、 y 两个词较接近，则 $p(w|x)$ 与 $p(w|y)$ 较接近。这是在 x 、 y 是符号表示时是难以建模到的！

大纲

■ 引言

■ 词表示

■ 神经概率语言模型 (NPLM)

■ 总结





■神经概率语言模型(NPLM)

■Bengio 2001, 2003模型

■A Neural Probabilistic Language Model, JMLR, 2003

■Bengio 2003 模型贡献

■得到基于分布式词表示的语言模型

■高阶(6), 无需平滑

■得到分布式词表示

■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 表变量, w_i 表具体词

■ v' 表转置, A_j 表 A 的第 j 行



■ NPLM: 两部分

■ 词表示

■ 牛 $\rightarrow C(\text{牛}) = (x_{11}, x_{12}, \dots, x_{1m})$

■ 马 $\rightarrow C(\text{马}) = (x_{21}, x_{22}, \dots, x_{2m})$

■ 网 $\rightarrow C(\text{网}) = (x_{31}, x_{32}, \dots, x_{3m})$

■ ...

■ 语言模型

■ $p(w|w_1, w_2, \dots, w_n) \rightarrow p(c(w)|c(w_1), c(w_2), \dots, c(w_n))$

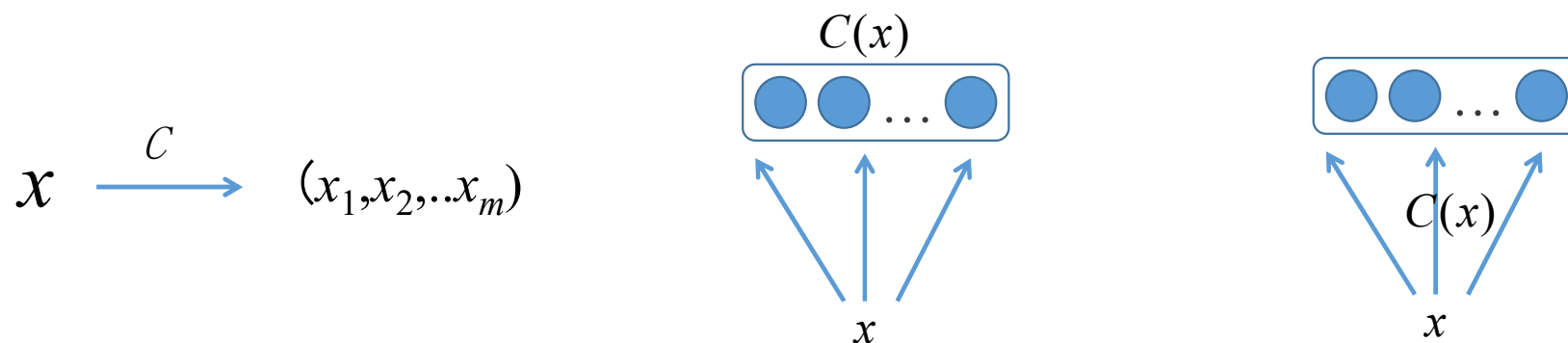
■ 1. 词表示:

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

■ 实现方式: 词表映射

■ 查表映射 C : 将任意词映射为一个 m 维向量

■ 或者说是一个 2 层的神经网络





■ 2. 语言模型 $p(c(x_t) | c(x_{t-1}), \dots, c(x_{t-(n-1)}))$

■ 建模连续函数:

$$\blacksquare c(x_{t-1}), \dots, c(x_{t-(n-1)}) \rightarrow p(c(x_t) | c(x_{t-1}), \dots, c(x_{t-(n-1)}))$$

■ 神经网络模型

■ 训练数据

■ 损失函数



■神经网络模型

■模型目标：训练一个映射 g 建模 n 元语言模型：

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

■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：

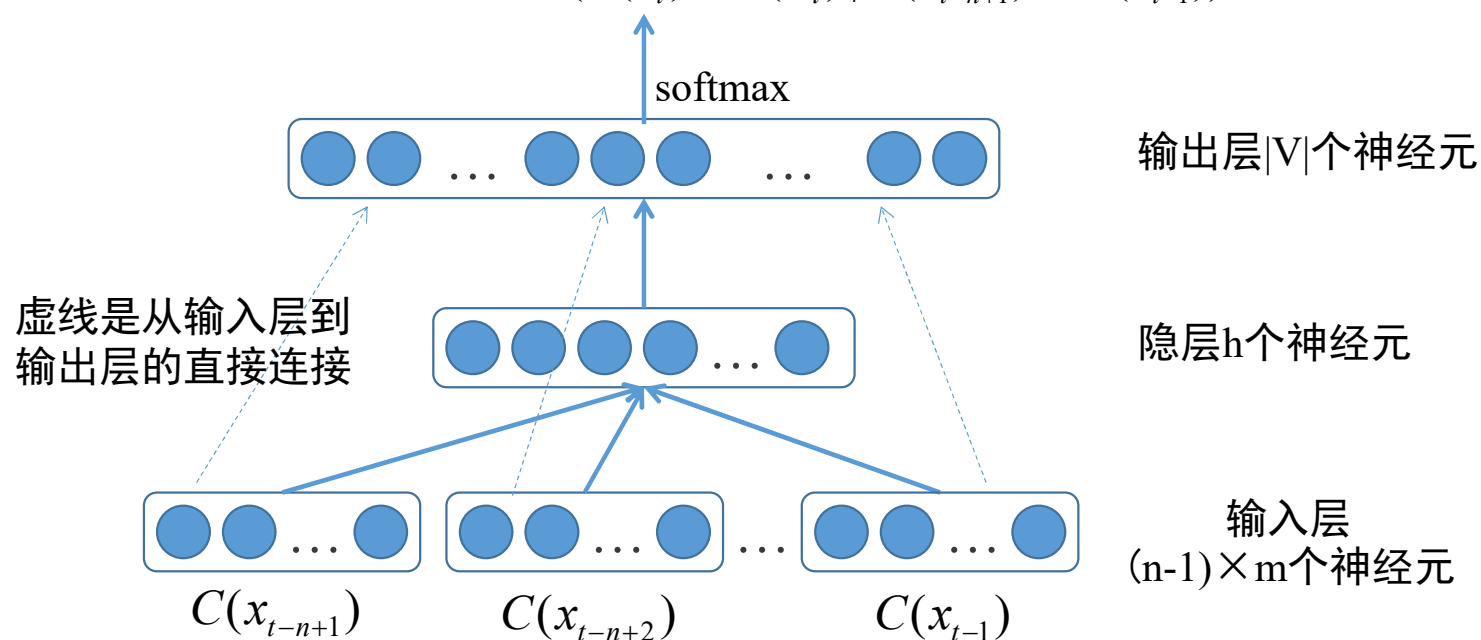
$$PP(C(x_1), \dots, C(x_T)) = P(C(x_1), \dots, C(x_T))^{-\frac{1}{T}} = \left(\prod_{t=1}^T P(C(x_t) | C(x_{t-1}), \dots, C(x_{t-n+1})) \right)^{-\frac{1}{T}}$$

■即极大化：

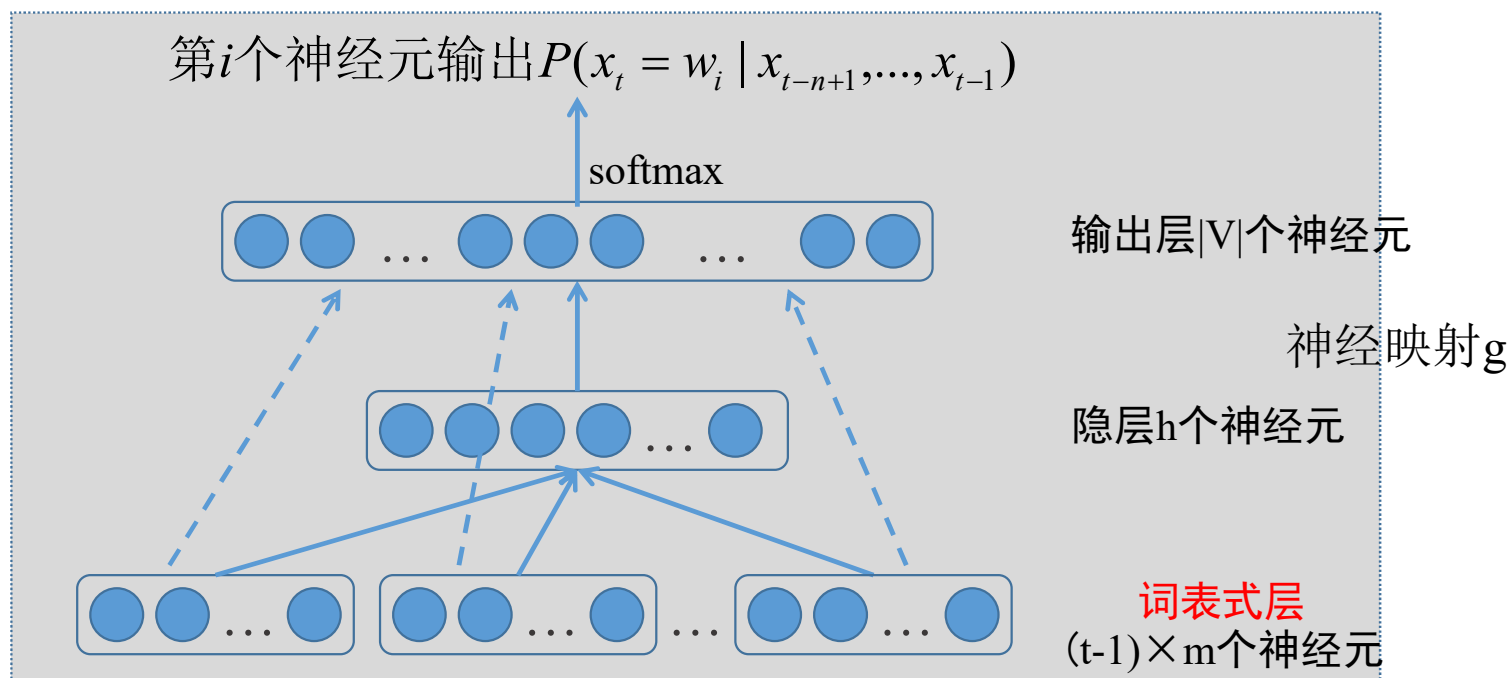
$$L = \frac{1}{T} \sum_{t=1}^T \log P(C(x_t) | C(x_{t-1}), \dots, C(x_{t-n+1})) = \frac{1}{T} \sum_{t=1}^T \log g(C(x_t), C(x_{t-1}), \dots, C(x_{t-n+1}); \omega)$$

神经网络模型结构

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



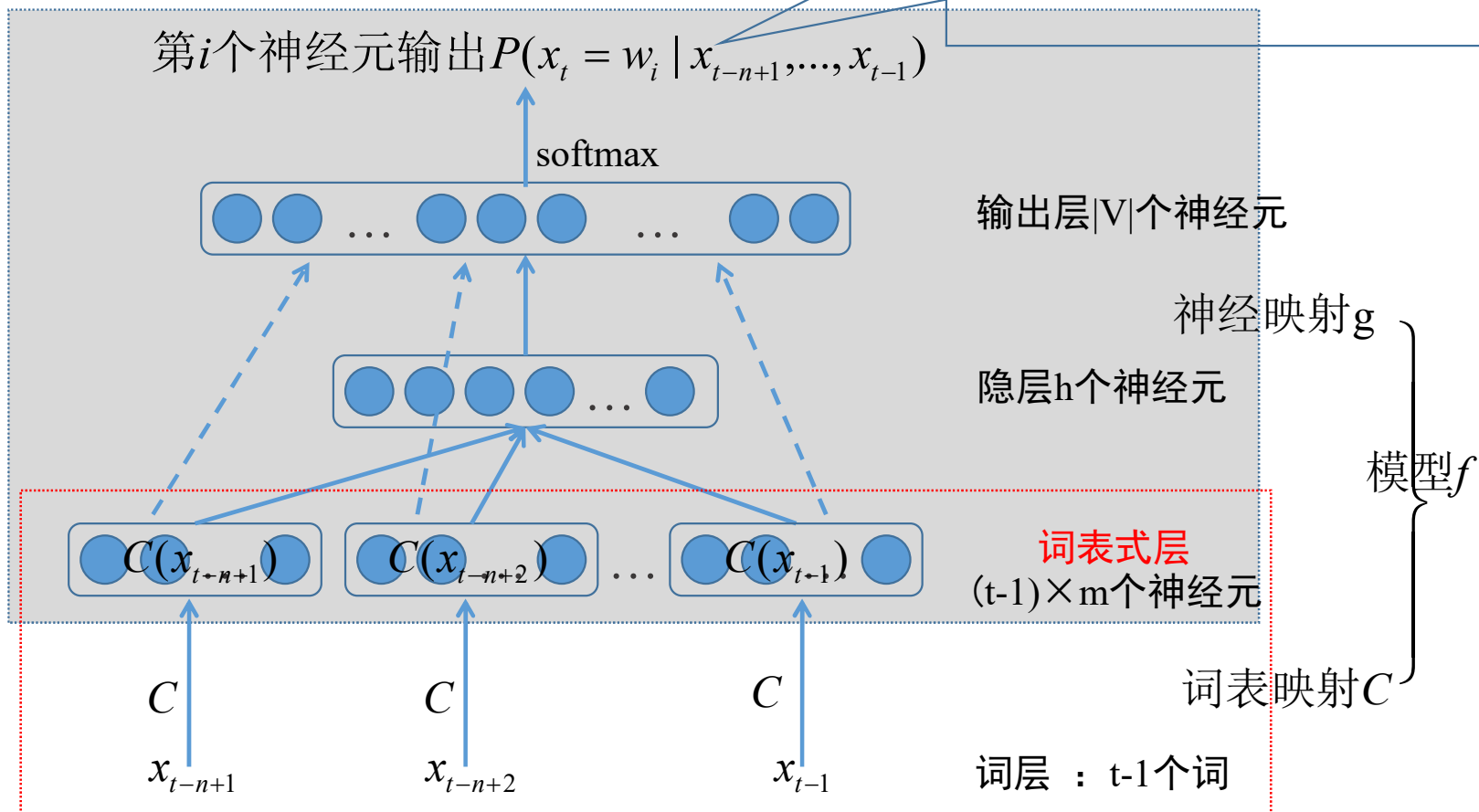
■合并词表映射与神经网络模型-1



合并词表映射与神经网络模型-2

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 元语法的 $n-1$ 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：每一个词都有表示， $C=|V| \times m$ 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： $y=b+W*C(x)+U\tanh(d+HC(x))$

■ 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



■混合模型

- 神经网络模型+插值trigram

■跨层连接

■算法的并行执行

- ...

■实验语料

	Brown语料	AP新闻
训练语料规模	1,181,041词的前800000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 $\varepsilon_0=10^{-3}$, 之后衰减, 按 $\varepsilon_t=\varepsilon_0/(1+rt)$	
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 的连接	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





■从上图得到的一些结论：

- 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)：
 - 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。



■ M&H 2007: Log-bilinear Language (LBL) model

Model type	Context size	Model test score	Mixture test score
Log-bilinear	5	117.0	97.3
Log-bilinear	10	<u>107.8</u>	<u>92.1</u>
Back-off KN3	2	129.8	
Back-off KN5	4	123.2↓	
Back-off KN6	5	123.5↑	
Back-off KN9	8	124.6↑	

■ Mnih, A., & Teh, Y. W. (2012)

Model	Context size	Latent Dim	Test PP
3-Gram	2		130.8
4-Gram	3		122.1
5-Gram	4		121.5
6-Gram	5		121.7
LBL	2	100	145.5
LBL	3	100	135.6
LBL	5	100	129.8
LBL	10	100	124.0
LBL	10*2	100	<u>38.6</u>
LBL	5	200	123.6
LBL	10	200	117.7
LBL	10*2	200	<u>33.6</u>
LBL	10	300	116.4



■神经概率语言模型(NPLM)

- Bengio 2001, 2003: NN-LM

 - A Neural Probabilistic Language Model, JMLR, 2003

 - M&H 2007, 2008: log-bilinear language (LBL) model

- 基于NN的问题: 固定输入长度

- Mikolov 2010, 2013: RNN-LM

 - Recurrent neural network based language model

RNN LM

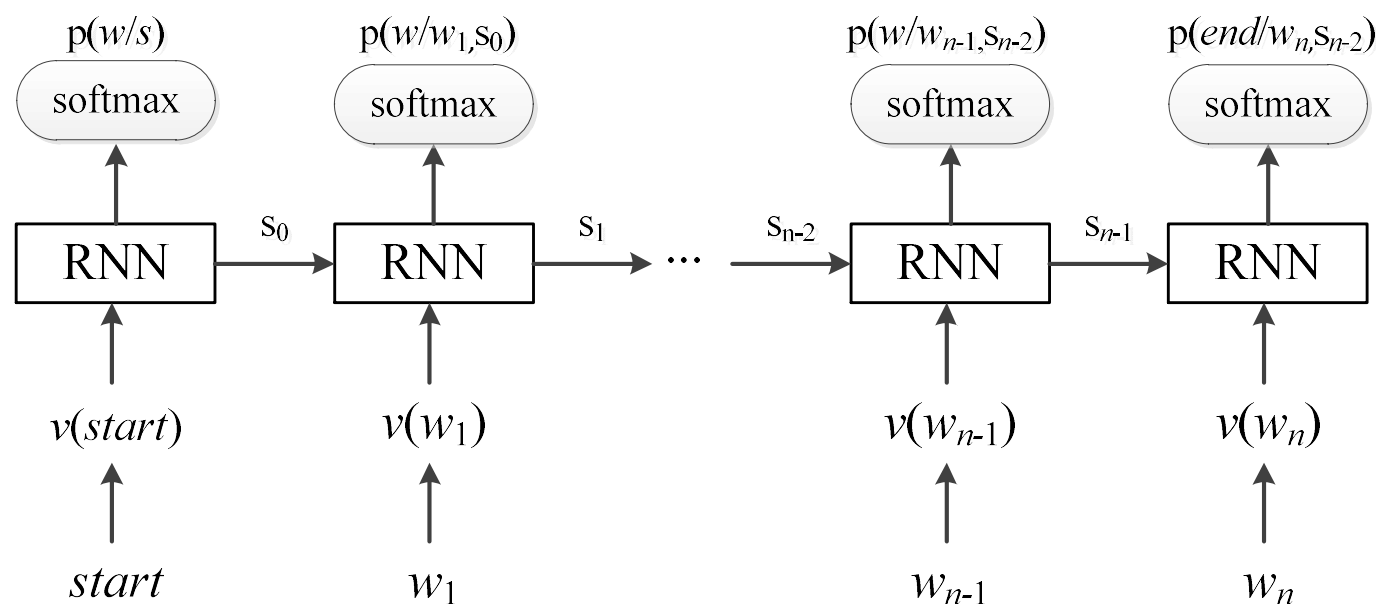
■ Mikolov, et al. Recurrent neural network based language model, Interspeech 2010

GT: $p(w_1/s)=1$,
other $w_i=0$

GT: if $w=w_i$, $p(w/w_{i-1})=1$, others=0

■ 不固定输入

■ 所以 N?





RNN LM

- 随训练数据规模增加时的情况
- 90：隐层节点数
- 2：合并低频词的词阈值

Model	#word	PPL
KN5 LM	200K	336
KN5 LM + RNN 90/2		271
KN5 LM	1M (0.8 6hours)	287
KN5 LM + RNN 90/2		225
KN5 LM	6.4M	221
KN5 LM + RNN 90/2		156



■神经语言模型的近期发展-1:

■从基于词单元的LM到基于其他不同语言单元:

- 英语: 基于character间关系
- 汉语: 基于字间关系
- 基于其他语言单元



■神经语言模型的近期发展-2

■从RNN LM到 **attention RNN LM**

■Mei2017, Salton 2017: attention RNN-LM

■从attention RNN LM到纯**self-attention LM(Transformer)**

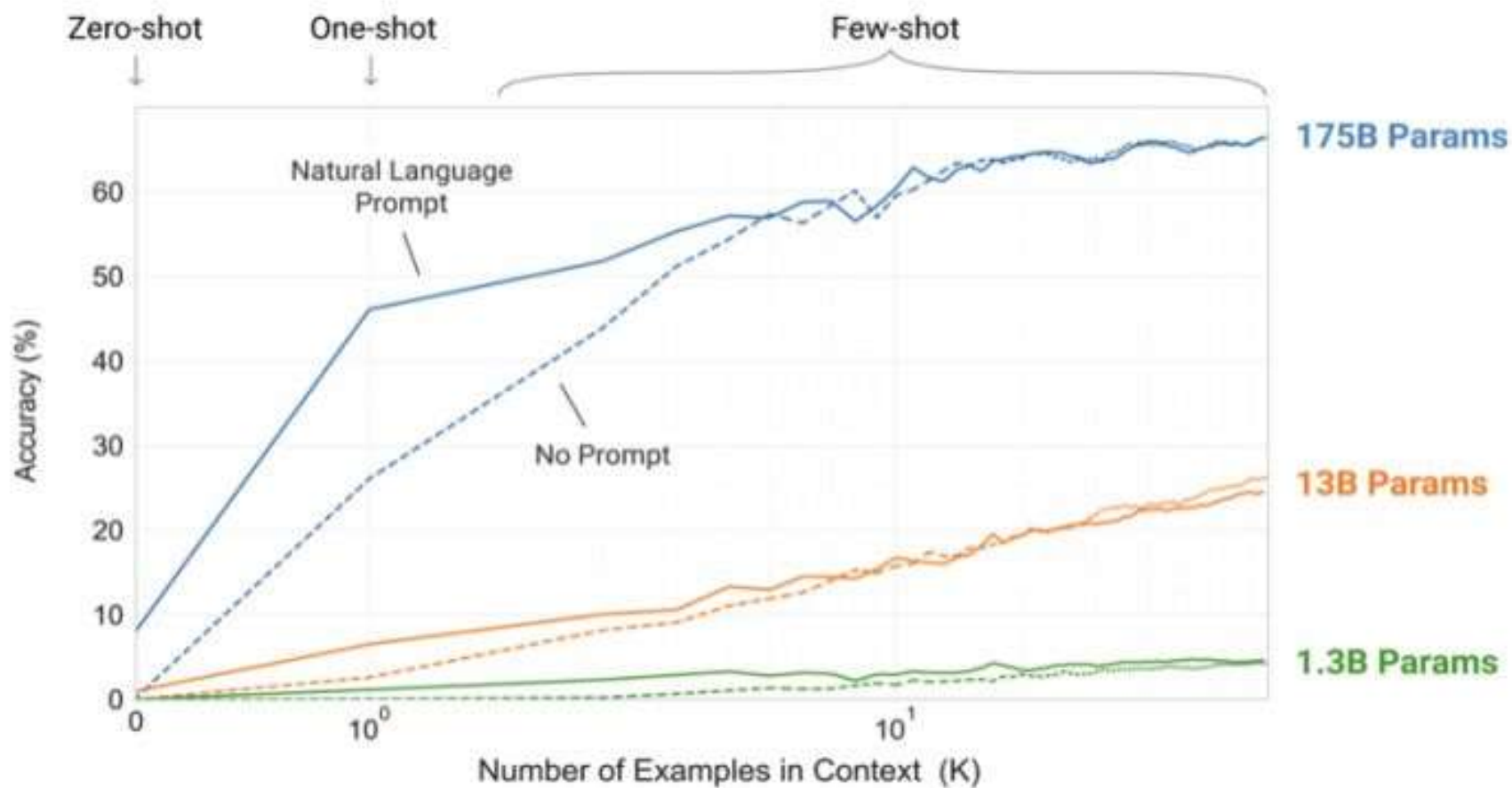
■GPT: Radford2018:improving language understanding by generative pre-training

■GPT2: Radford 2019:Language Models are Unsupervised Multitask Learners (48 层 Transformer Decoder, 15亿参数)

■GPT3: 2020 (96 层 Transformer Decoder, 1750亿参数(10^{11}))

■<https://beta.openai.com/>

数据稀疏情况下



大纲

■ 引言

■ 词表示

■ 神经概率语言模型 (NPLM)

■ 总结





■总结

- 从基于MLE的LM到基于神经网络的LM

- 表示上：从符号词表示到向量词表示

 - 关键

- 估计上：从基于统计的估计到基于预测的估计

 - 引入更多变化



思考与讨论

- N-gram语言模型建模了什么知识？
- 这种知识有什么用处？
- N-gram语言模型建模的知识有什么问题？
- 扩大N的其他途径：组块
- 灵活选择N的途径：递归



Thank You !