

# 深度学习-自然语言处理

史春奇

2017 年

# 目录

<b>1 自然语言处理简介</b>	<b>6</b>
1.1 Seq2Seq 模型	8
1.2 RNN Encoder-Decoder 模型	12
1.3 Encoder-Decoder 模型	15
1.4 NNLM 模型	18
1.5 Log-Linear 模型	20
1.6 嵌入词模型	22
1.7 向量空间模型	30
1.8 词袋 BoW 模型	37
1.9 N-Gram 模型	40
<b>2 数学基础</b>	<b>46</b>
2.1 概率统计	46
2.2 频率派	48

2.3 Fisher 派 . . . . .	52
2.4 贝叶斯派 . . . . .	59
2.5 K 分类最大相互熵推导 Softmax . . . . .	65
<b>3 词嵌入模型 . . . . .</b>	<b>70</b>
3.1 词-文上下文 Word-Document . . . . .	70
3.2 词上下文 Word-Neighboring Word . . . . .	78
3.3 Softmax 近似计算 . . . . .	116
<b>4 语言模型 . . . . .</b>	<b>123</b>
4.1 从概率模型到神经网络模型 . . . . .	123
4.2 从神经网络模型到 RNN 模型 . . . . .	140
4.3 从 RNN 模型到 Encoder-Decoder 模型 . . . . .	147
4.4 Attention 模型 . . . . .	164
4.5 Seq2Seq 模型：端到端 . . . . .	184
4.6 翻译输出 . . . . .	192
4.7 替代和增强 RNN 模型 . . . . .	194

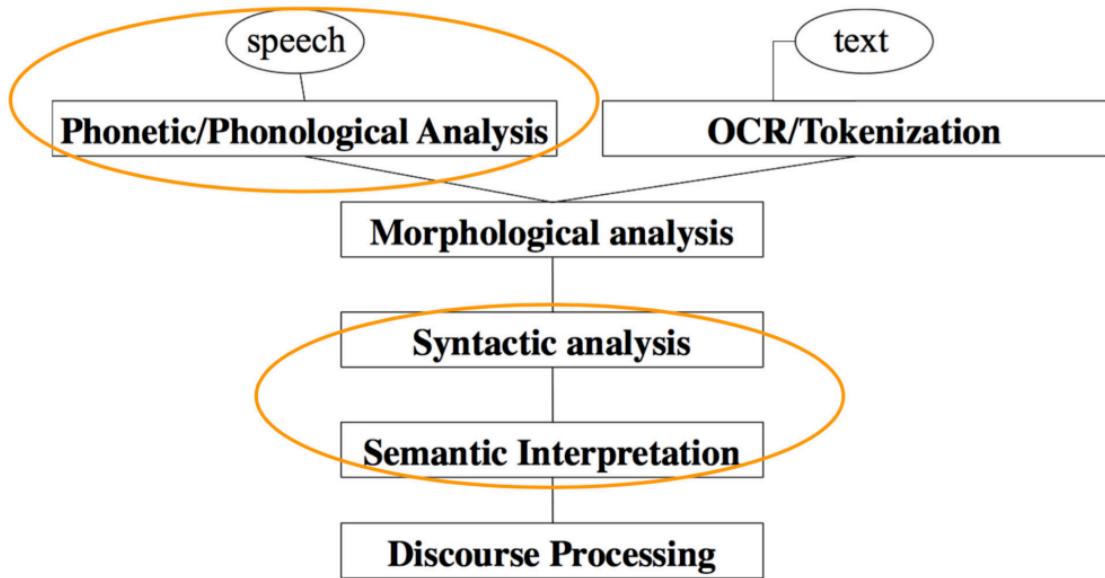
4.8 其他应用 . . . . .	207
--------------------	-----

## 第三部分-自然语言处理

- 深度自然语言处理的由来
  - 谁发明了神经网络语言模型？
  - 谁最早应用 RNN 到语言模型？
  - 词嵌入是深度模型么？
- 词嵌入的理解
  - 同样是无监督，词嵌入为什么优于主题模型？
  - 词类比是如何可以做到的？
  - 上下文和上文有什么区别？
- 神经网络机器翻译
  - 编码器和解码器为什么能做机器翻译？
  - 为什么要使用双向 RNN？
  - 什么是软（硬）注意力机制？

# 1 自然语言处理简介

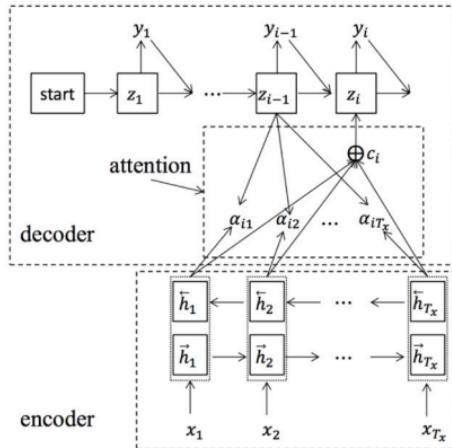
- NLP 范围



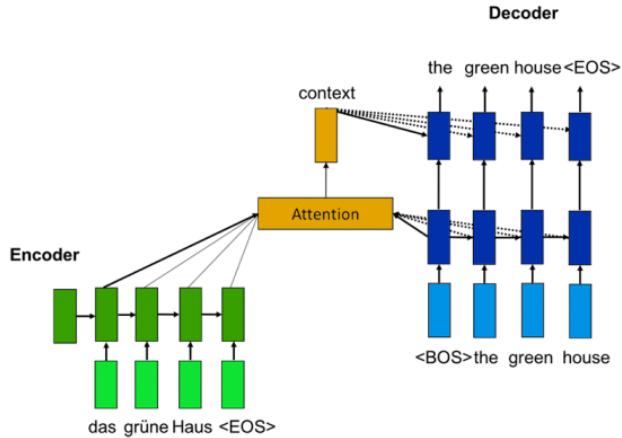
- NLP 范围
  - 1. 情感分析 Sentiment Analysis
  - 2. 问答 Question Answering
  - 3. 机器翻译 Machine Translation
- 最有代表性的自然语言处理

## 1.1 Seq2Seq 模型

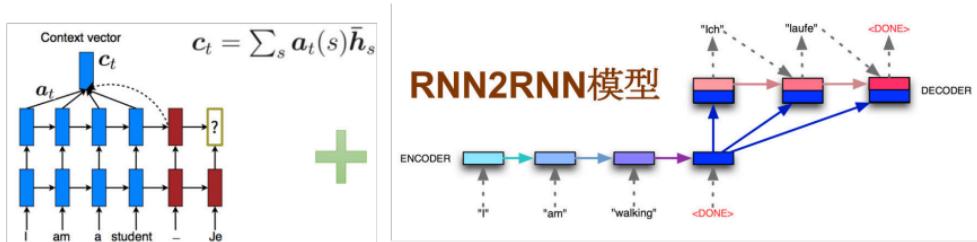
- Seq2Seq 机器翻译模型



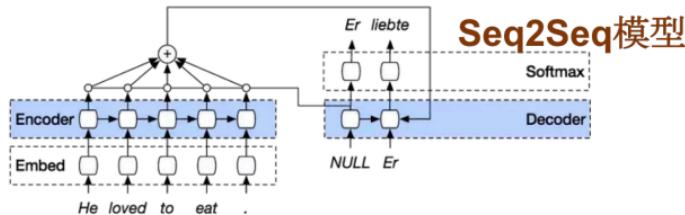
- Seq2Seq 机器翻译模型



- Seq2Seq 机器翻译模型: RNN + Attention

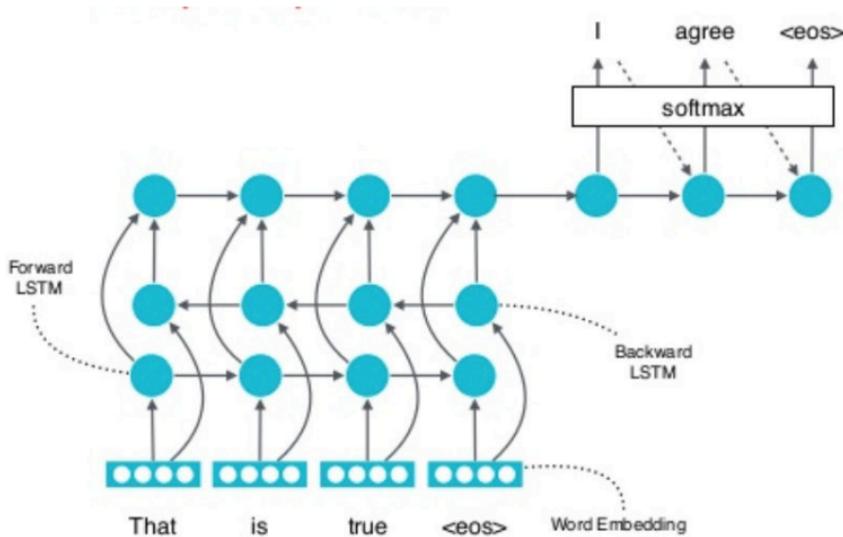


Attention模型



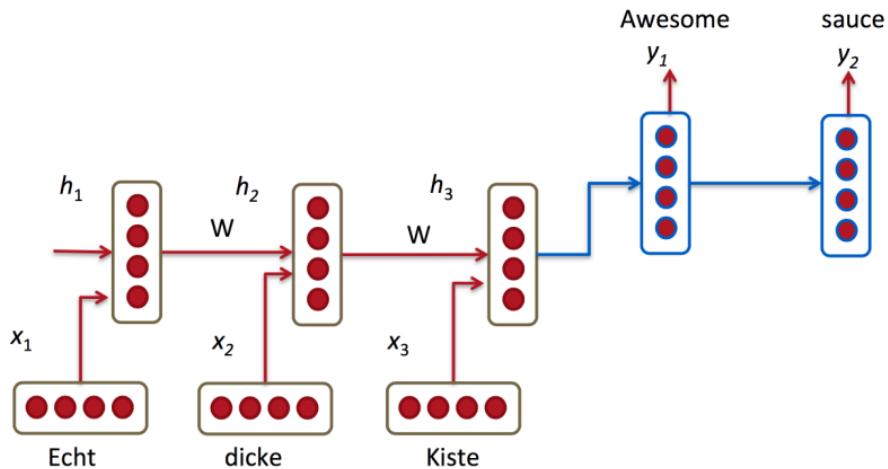
Seq2Seq模型

- Seq2Seq: Bidirectional RNN + Attention

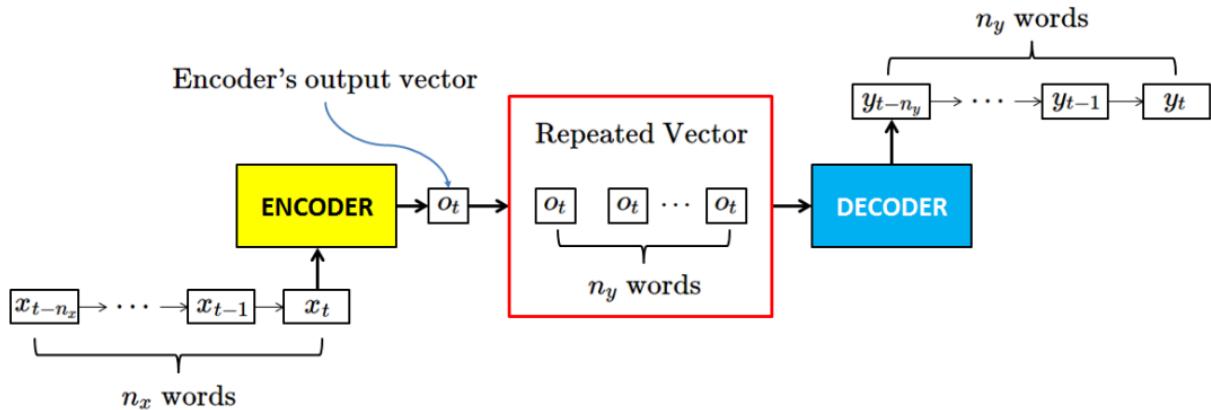


## 1.2 RNN Encoder-Decoder 模型

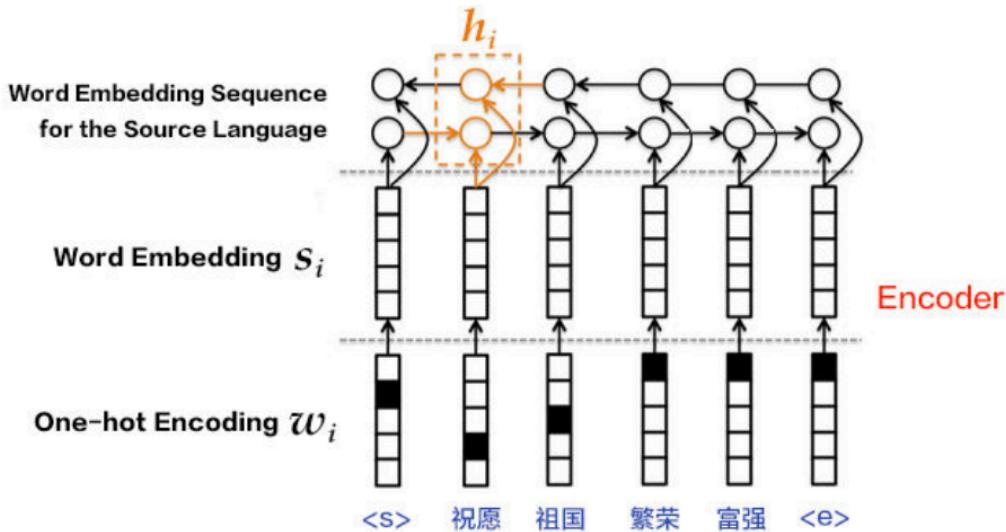
- RNN 模型



- RNN Encoder-Decoder

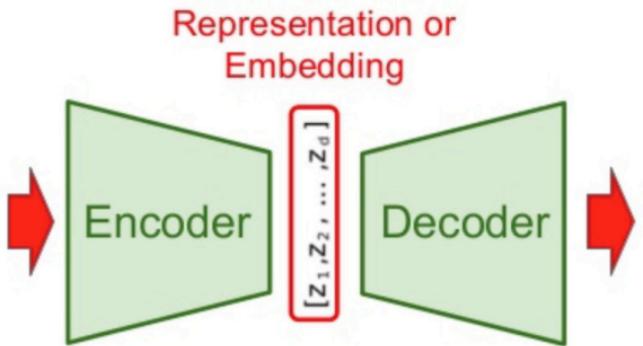


- Bidirectional RNN Encoder



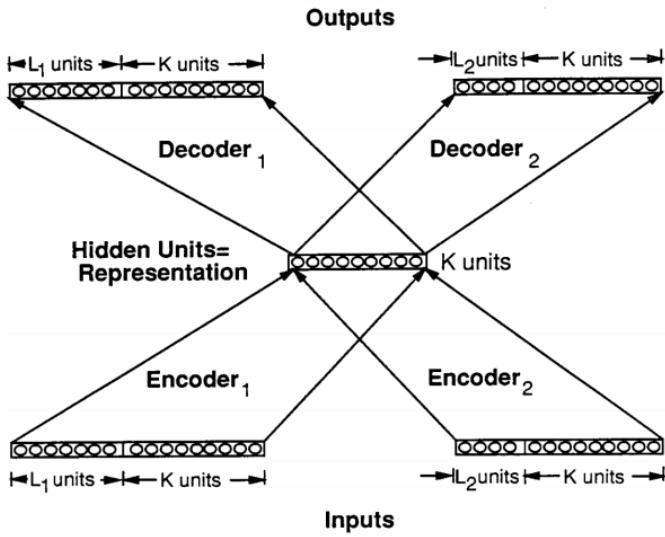
### 1.3 Encoder-Decoder 模型

Economic growth has slowed down in recent years.

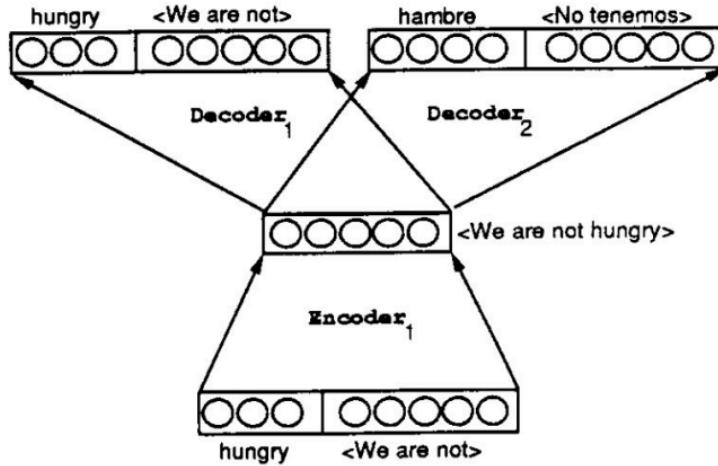
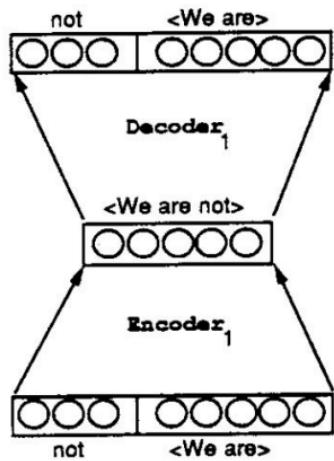


La croissance économique a ralenti ces dernières années.

- Dual-Ported Encoder-Decoder

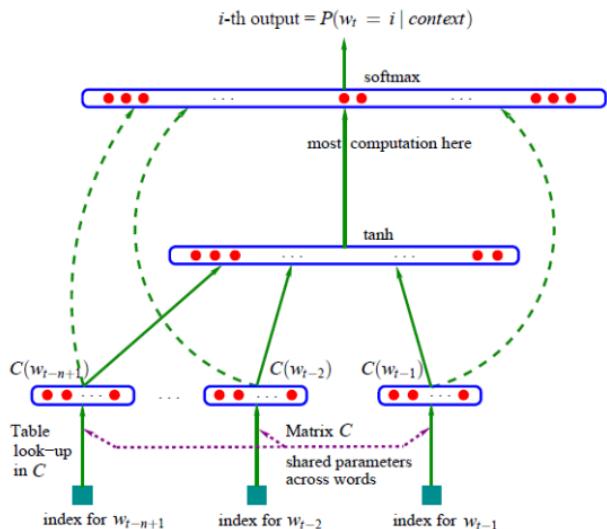


- Contextual Encoder-Decoder



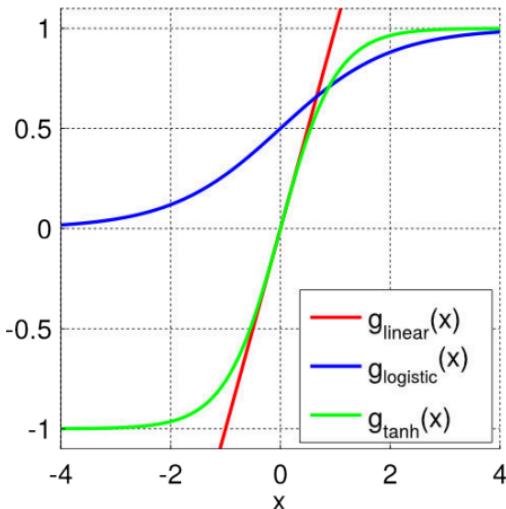
## 1.4 NNLm 模型

- 神经网络语言模型 (Bengio)

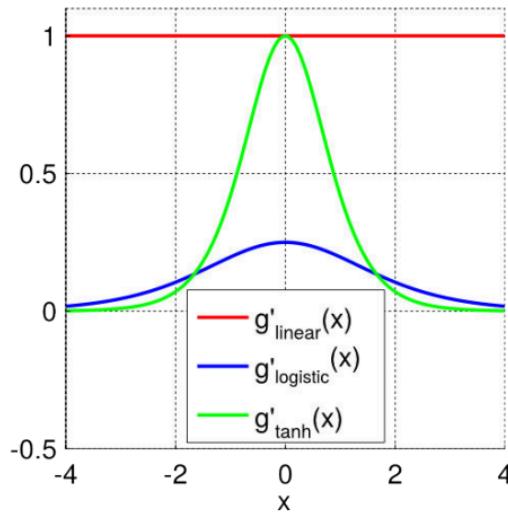


- Tanh 在 0 附近近似  $y=x$

Some Common Activation Functions



Activation Function Derivatives



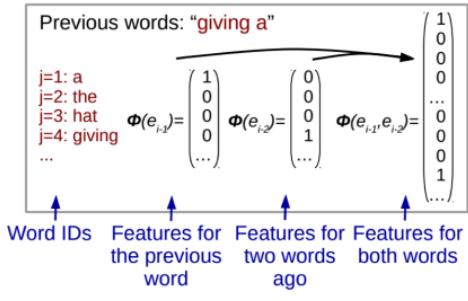
## 1.5 Log-Linear 模型

- Log-Linear 模型

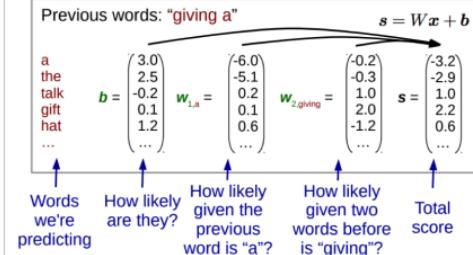
$$P(y \mid x, \mathbf{W}) = \frac{e^{\sum_k \mathbf{W}_k \phi_k(x, y)}}{\sum_{y' \in \mathcal{Y}} e^{\sum_k \mathbf{W}_k \phi_k(x, y')}}$$

概率  $p = \text{softmax}(s)$

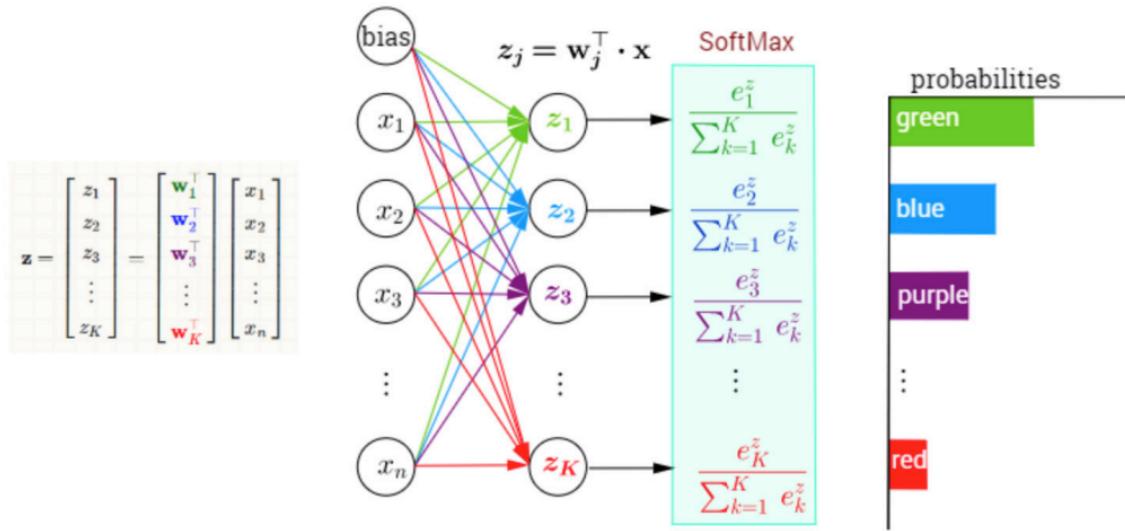
$$p_j = \frac{\exp(s_j)}{\sum_j \exp(s_j)}$$



限制

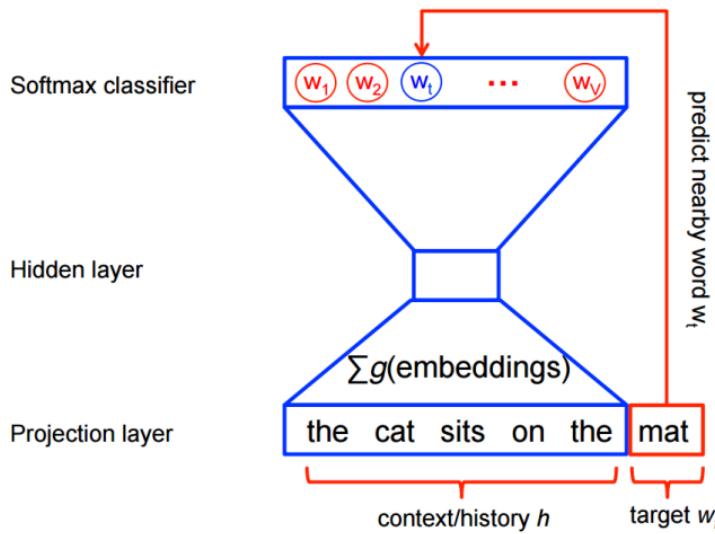


- Softmax 神经网络模型

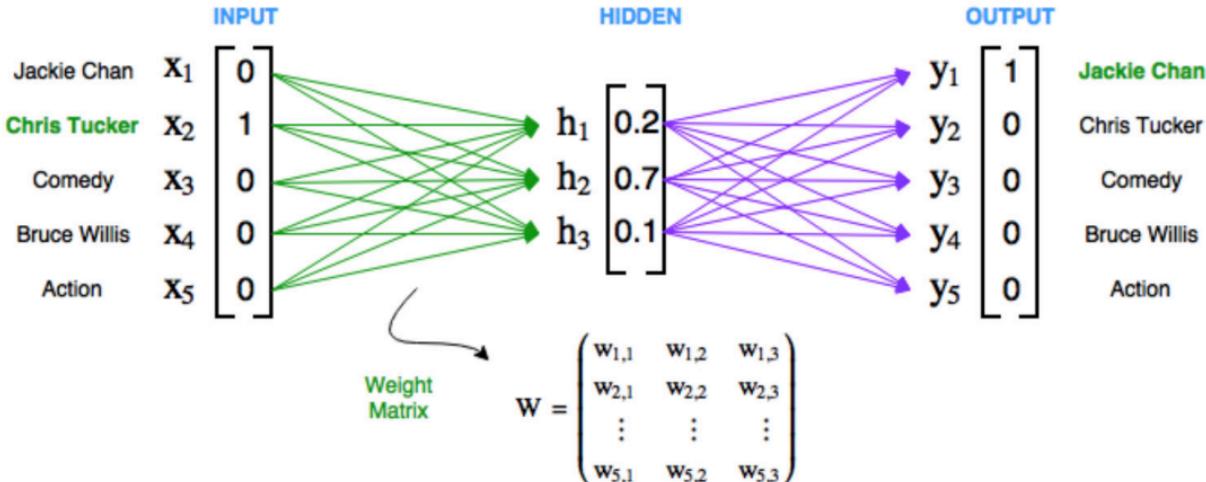


## 1.6 嵌入词模型

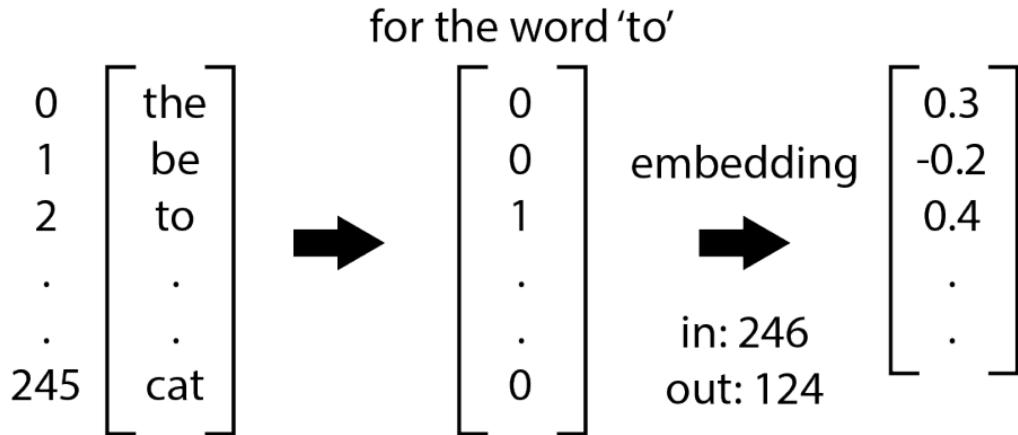
- Softmax 神经网络模型



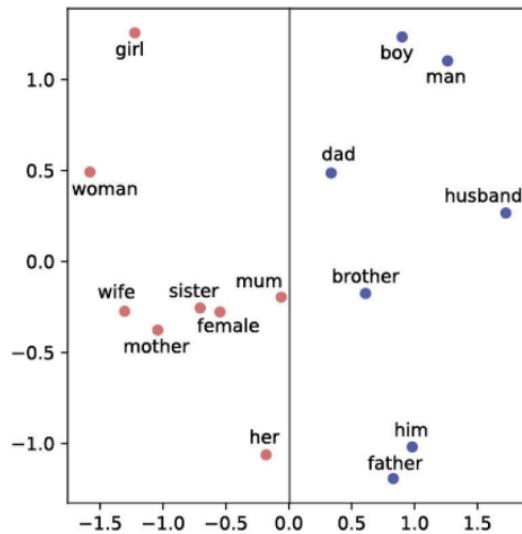
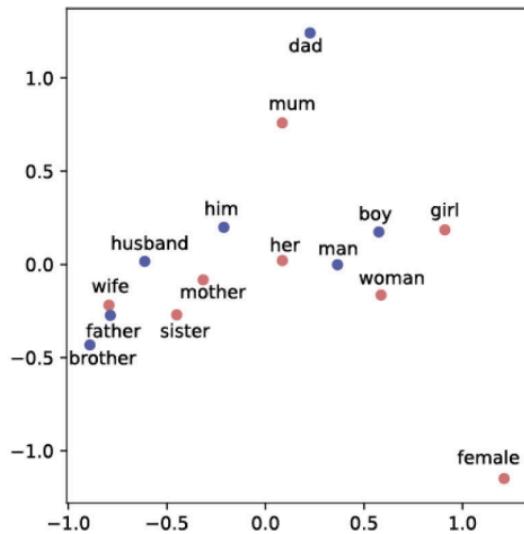
- 词嵌入模型



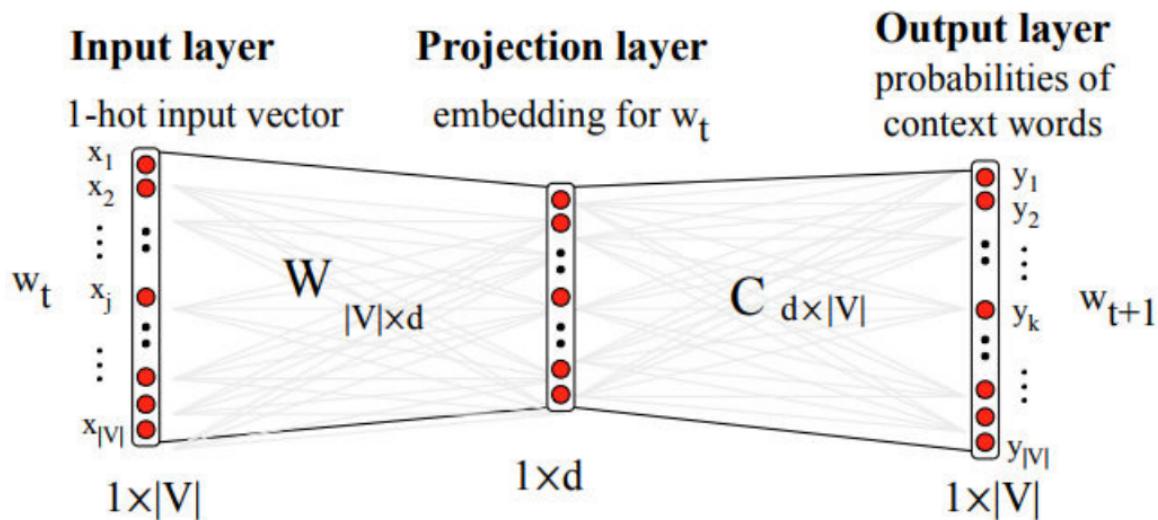
- 词嵌入模型



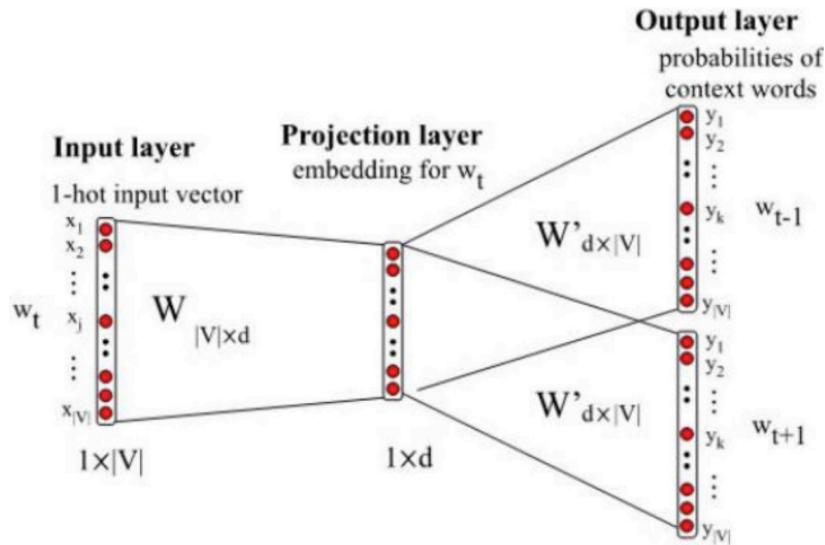
- 词嵌入模型例子



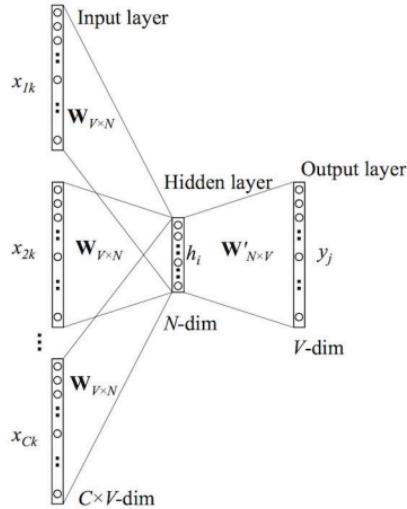
- 词嵌入模型例子



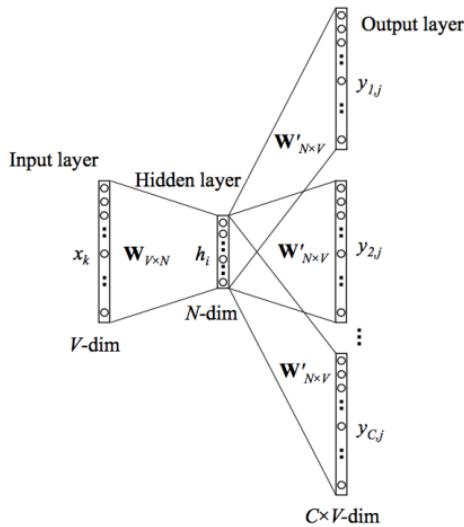
- 词嵌入模型例子



- 连续词袋 CBoW 模型

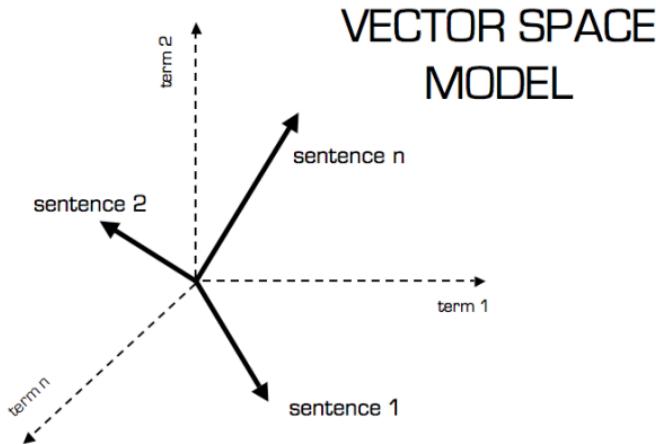


- Skip-Gram 模型



## 1.7 向量空间模型

- 向量空间 Vector Space Model: 句子



- TF-IDF

$$w_{x,y} = \text{tf}_{x,y} \times \log \left( \frac{N}{df_x} \right)$$

## TF-IDF

Term  $x$  within document  $y$

$\text{tf}_{x,y}$  = frequency of  $x$  in  $y$

$df_x$  = number of documents containing  $x$

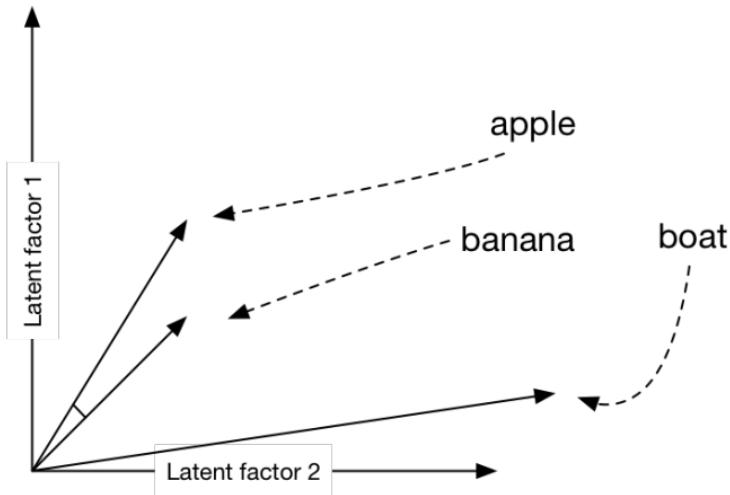
$N$  = total number of documents

	Document 1	Document 2	Document 3	Document 4	Document 5	Document 6	Document 7	Document 8
Term(s) 1	10	0	1	0	0	0	2	
Term(s) 2	0	2	0	0	0	18	0	2
Term(s) 3	0	0	0	0	0	0	0	2
Term(s) 4	6	0	0	4	6	0	0	0
Term(s) 5	0	0	0	0	0	0	0	2
Term(s) 6	0	0	1	0	0	1	0	0
Term(s) 7	0	1	8	0	0	0	0	0
Term(s) 8	0	0	0	0	0	3	0	0

↑  
Document Vector

Word Vector  
(Passage Vector)

- 向量空间 Vector Space Model: 词



- 共现矩阵 Co-occurrence Matrix

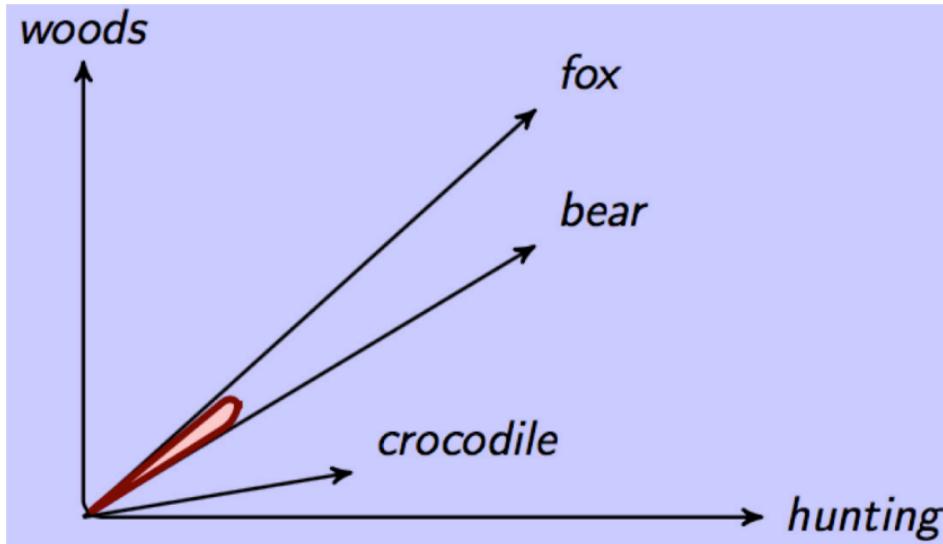
The diagram illustrates the construction of a co-occurrence matrix. Three sets of words are shown in colored boxes:

- Top row (blue border): car engine hood tires truck trunk
- Middle row (red border): car emissions hood make model trunk
- Bottom row (orange border): Chomsky corpus noun parsing tagging wonderful

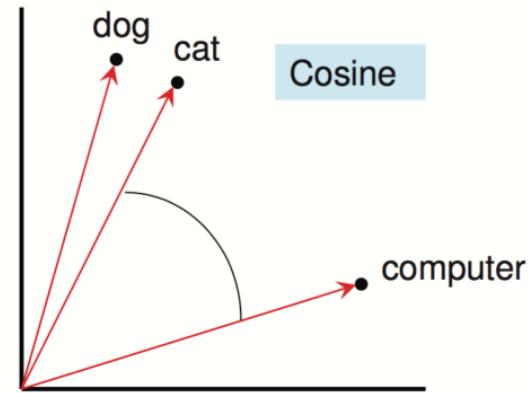
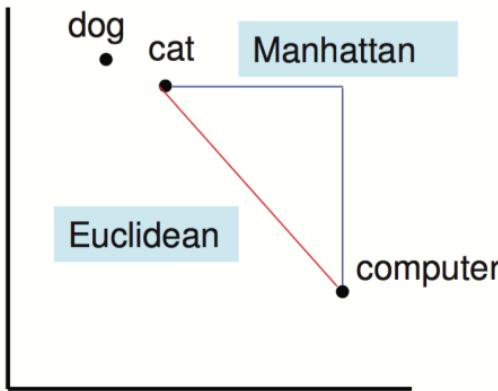
A green arrow points downwards from these words to a 3x15 grid labeled "共现矩阵" (Co-occurrence Matrix) in red.

	car	Chomsky	corpus	emissions	engine	hood	make	model	noun	parsing	tagging	tires	truck	trunk	wonderful
car	0	0	0	0	1	1	0	0	0	0	0	1	1	1	0
hood	1	0	0	1	0	0	1	1	0	0	0	0	0	1	0
Chomsky	0	0	1	0	0	0	0	0	1	1	1	0	0	0	1

- Cosine 距离



- 其他常用距离：欧式距离和曼哈顿距离



- 距离计算公式

The *cosine* of the angle between two vectors  $\mathbf{x}$  and  $\mathbf{y}$  is:

$$\cos(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x} \cdot \mathbf{y}}{||\mathbf{x}|| \cdot ||\mathbf{y}||} = \frac{\sum_{i=1}^n x_i y_i}{\sqrt{\sum_{i=1}^n x_i^2} \sqrt{\sum_{i=1}^n y_i^2}}$$

The *Euclidean distance* of two vectors  $\mathbf{x}$  and  $\mathbf{y}$  is:

$$||\mathbf{x} - \mathbf{y}|| = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

## 1.8 词袋 BoW 模型

- BoW 模型：基于词典序列

"Represent each document which the bag of words it contains"

d1 : Mary loves Movies, Cinema and Art

Class 1 : Arts

d2 : John went to the Football game

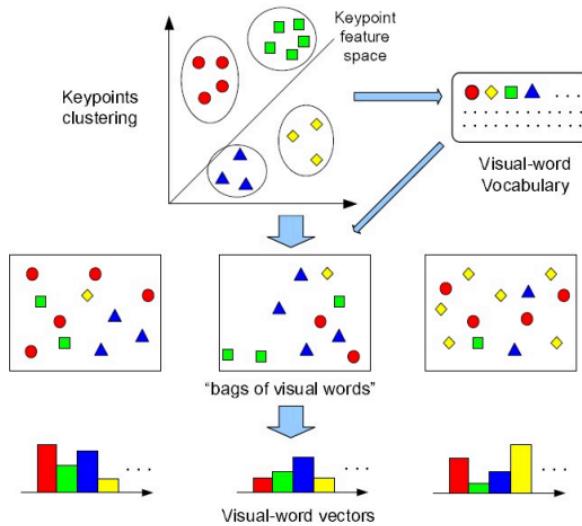
Class 2 : Sports

d3 : Robert went for the Movie Delicatessen

Class : Arts

	Mary	Loves	Movies	Cinema	Art	John	Went	to	the	Delicatessen	Robert	Football	Game	and	for
d1	1	1	1	1	1										1
d2						1	1	1	1			1	1		
d3			1				1		1		1				1

- BoW 模型

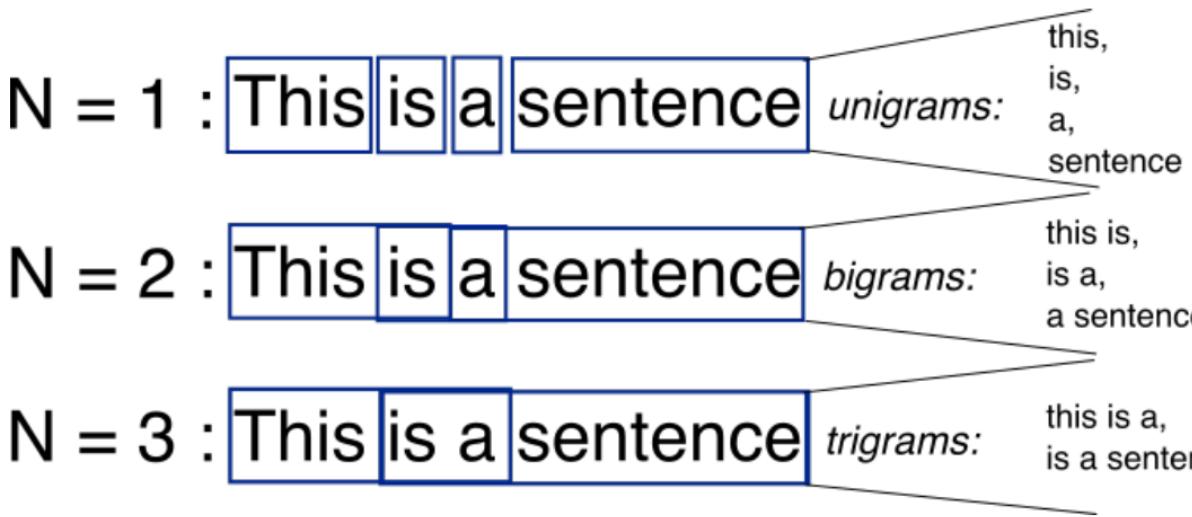


- 1-Hot-Encoding

Original data:		One-hot encoding format:					
id	Color	id	White	Red	Black	Purple	Gold
1	White	1	1	0	0	0	0
2	Red	2	0	1	0	0	0
3	Black	3	0	0	1	0	0
4	Purple	4	0	0	0	1	0
5	Gold	5	0	0	0	0	1

## 1.9 N-Gram 模型

- N-Gram



- N-Gram

unigram



bigram



trigram



n-gram ( $n = 4$ )



- 全概率公式展开

$$P(Girl | Flower) = \frac{\sum count(Flower Girl)}{\sum count(Girl)}$$

$$P(W) = P(w_1, w_2, w_3, w_4, w_5 \dots w_n)$$

$$P(x_1, x_2, x_3, \dots, x_n) = P(x_1)P(x_2 | x_1)P(x_3 | x_1, x_2) \dots P(x_n | x_1, \dots, x_{n-1})$$

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

P("its water is so transparent") =

$$P(\text{its}) \times P(\text{water} | \text{its}) \times P(\text{is} | \text{its water})$$

$$\times P(\text{so} | \text{its water is}) \times P(\text{transparent} | \text{its water is so})$$

- 马尔可夫近似

1阶Markov假设 (bigram) :

$$P(\text{the} \mid \text{its water is so transparent that}) \approx P(\text{the} \mid \text{that})$$

2阶Markov假设 (trigram) :

$$P(\text{the} \mid \text{its water is so transparent that}) \approx P(\text{the} \mid \text{transparent that})$$

$$P(\text{Girl} \mid \text{Flower}) = \frac{\Sigma \text{count(Flower Girl)}}{\Sigma \text{count(Girl)}}$$

- 概率估算问题

$$\begin{aligned} P(|E|=3, e_1="she", e_2="went", e_3="home") = \\ & P(e_1="she") \\ & * P(e_2="went" | e_1="she") \\ & * P(e_3="home" | e_1="she", e_2="went") \\ & * P(e_4="</s>" | e_1="she", e_2="went", e_3="home") \end{aligned}$$

$$P(E) = \prod_{t=1}^{T+1} P(e_t | e_1^{t-1})$$

高阶 N-Gram 数据稀疏，不存在问题！



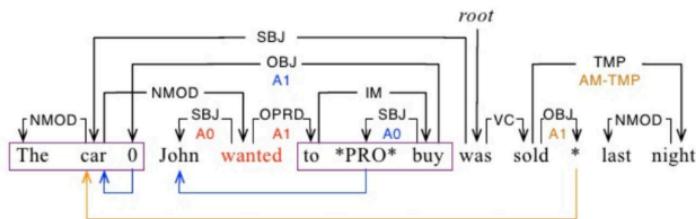
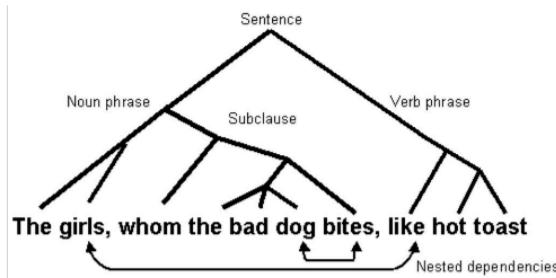
### Smoothing 技术：线性插值低阶N-Gram

$$P(e_t | e_{t-1}) = (1 - \alpha)P_{ML}(e_t | e_{t-1}) + \alpha P_{ML}(e_t)$$

Modified Kneser-Ney smoothing

$$P(e_t | e_{t-m+1}^{t-1}) = (1 - \alpha_m)P_{ML}(e_t | e_{t-m+1}^{t-1}) + \alpha_m P(e_t | e_{t-m+2}^{t-1})$$

- 长距离依赖问题

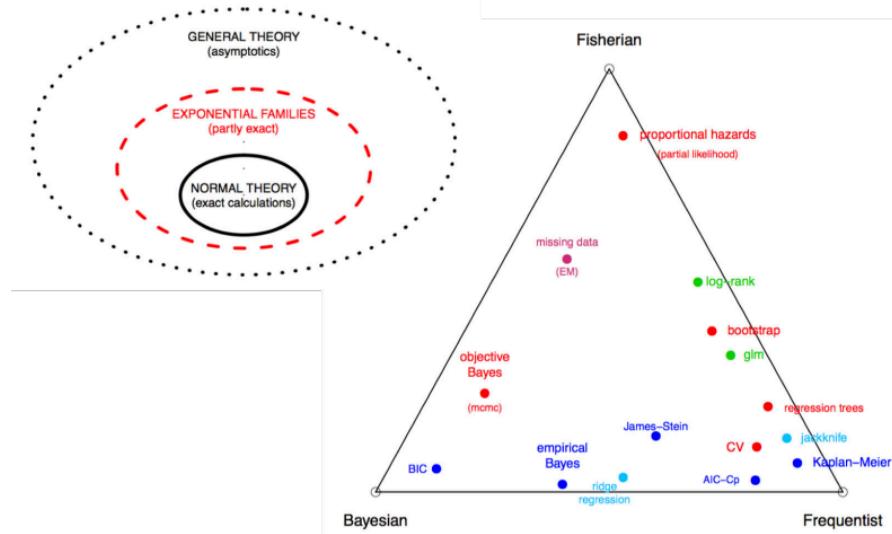


## 2 数学基础

### 2.1 概率统计

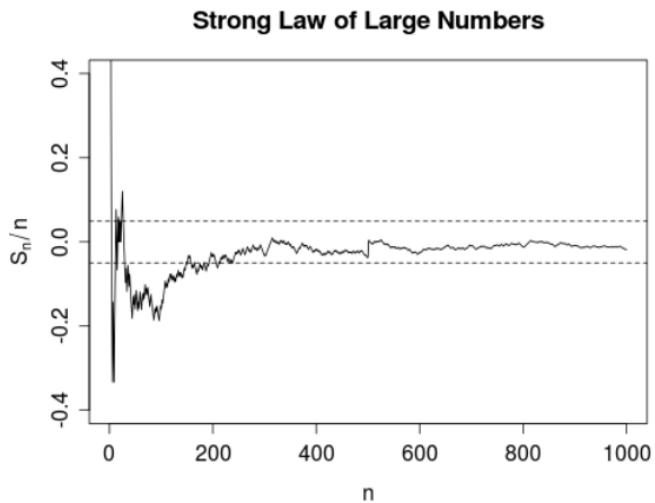
- 三大学派
  - 1. 频率派基于计数，通过概率极限建立概率分布
  - 2. Fisher 派基于经验，通过最大似然估计建立概率分布
  - 3. 贝叶斯派基于推理，通过最大熵和贝叶斯公式建立概率分布

- 分布、方法、学派

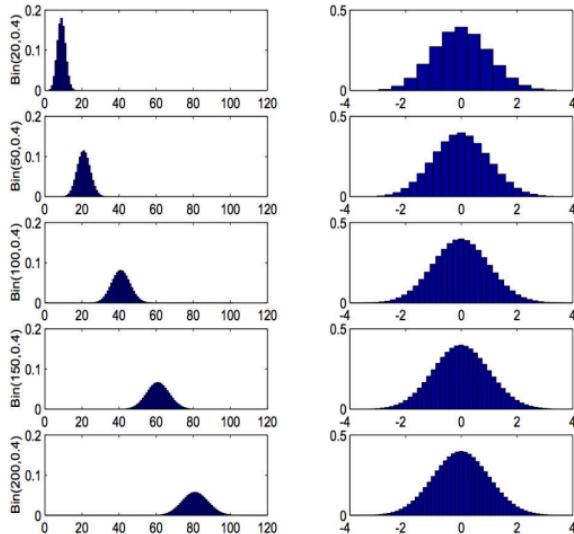


## 2.2 频率派

- 大数定理和按概率收敛



- 二项分布的大数极限是正态分布



- 二项分布的大数极限是正态分布

## 二项分布到正态分布

$$b(n, p, k) = \mathbb{P}\{Bin(n, p) = k\}$$
$$b(i) = b(n, \frac{1}{2}, i) = \binom{n}{i} (\frac{1}{2})^n$$

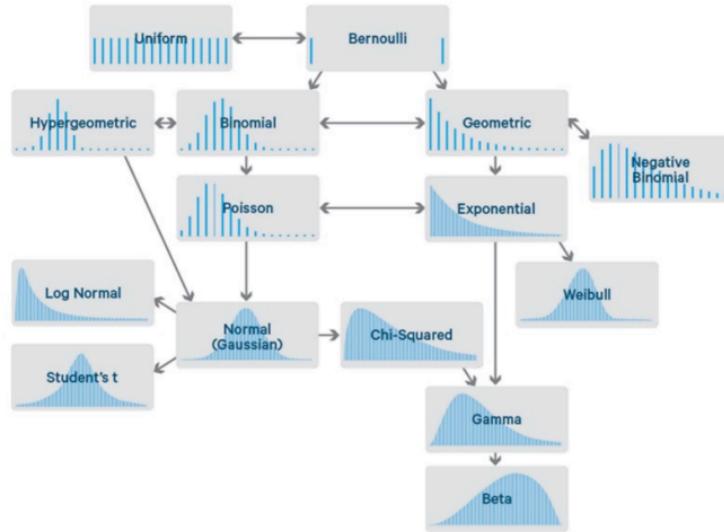
$$\frac{C_{n/2+1}^n}{C_{n/2}^n} = e^{\frac{-2i^2}{n}}$$

$$\frac{C_{n/2}^n}{2^n} \cong \frac{2}{\sqrt{2\pi n}}$$

$$\frac{b(\frac{n}{2} + d)}{b(\frac{n}{2})} \sim e^{-\frac{2d^2}{n}}$$

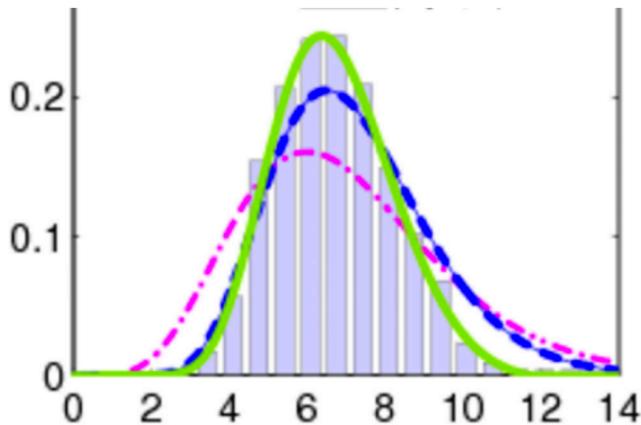
$$\sum_{-c\sqrt{n} \leq i \leq c\sqrt{n}} b\left(\frac{n}{2} + i\right)$$
$$\sim \sum_{-c\sqrt{n} \leq i \leq c\sqrt{n}} \frac{2}{\sqrt{2\pi n}} e^{-\frac{2i^2}{n}}$$
$$= \sum_{-2c \leq \frac{2i}{\sqrt{n}} \leq 2c} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(\frac{2i}{\sqrt{n}})^2} \frac{2}{\sqrt{n}}$$
$$\sim \int_{-2c}^{2c} \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx$$

- 正态分布到指数族分布



## 2.3 Fisher 派

- 最大似然估计和经验

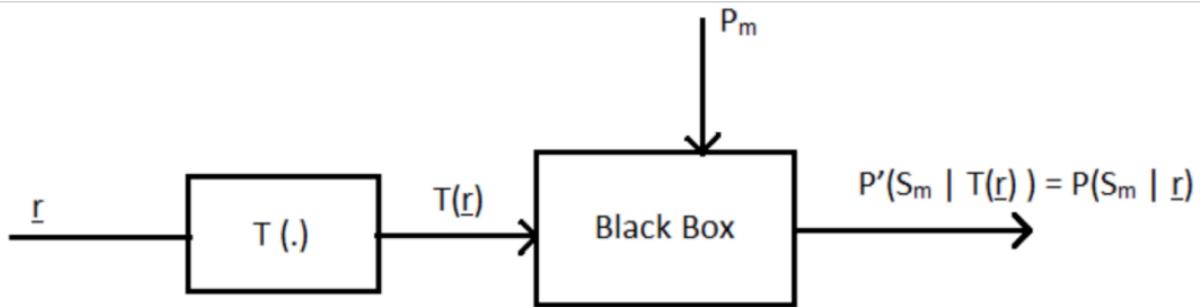


- 奥卡姆剃刀原理



- 充分统计量

$$L(x_1, x_2, \dots, x_n | \vartheta) = g(\vartheta, T(x_1, x_2, \dots, x_n)) \cdot h(x_1, x_2, \dots, x_n)$$



- Fisher-Neyman 分解

Fisher–Neyman factorization theorem

$$f_{\theta}(x) = h(x) g_{\theta}(T(x))$$

$$\begin{aligned} f_{X_1^n}(x_1^n) &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x_i - \theta)^2}{2\sigma^2}\right) \\ &= (2\pi\sigma^2)^{-\frac{n}{2}} \exp\left(-\sum_{i=1}^n \frac{(x_i - \theta)^2}{2\sigma^2}\right) \\ &= (2\pi\sigma^2)^{-\frac{n}{2}} \exp\left(-\sum_{i=1}^n \frac{((x_i - \bar{x}) - (\theta - \bar{x}))^2}{2\sigma^2}\right) \\ &= (2\pi\sigma^2)^{-\frac{n}{2}} \exp\left(-\frac{1}{2\sigma^2} \left(\sum_{i=1}^n (x_i - \bar{x})^2 + \sum_{i=1}^n (\theta - \bar{x})^2 - 2 \sum_{i=1}^n (x_i - \bar{x})(\theta - \bar{x})\right)\right) \\ &= (2\pi\sigma^2)^{-\frac{n}{2}} \exp\left(-\frac{1}{2\sigma^2} \left(\sum_{i=1}^n (x_i - \bar{x})^2 + n(\theta - \bar{x})^2\right)\right) \\ &= (2\pi\sigma^2)^{-\frac{n}{2}} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \bar{x})^2\right) \exp\left(-\frac{n}{2\sigma^2} (\theta - \bar{x})^2\right) \end{aligned}$$

- 最大似然估计到正态分布

$x_1, \dots, x_n$  n次观测值

$e_i = x_i - \theta$  n次观测误差

假定误差的密度分布函数为  $f(e)$

$$\frac{d}{dx}(e^x) = e^x$$

$$\frac{d}{dx}(a^x) = a^x \ln a$$

$$\frac{d}{dx}(\ln x) = \frac{1}{x}$$

$$\frac{d}{dx}(\log_a x) = \frac{1}{x \ln a}$$

$$L(\theta) = L(\theta; x_1, \dots, x_n)$$

$$= f(e_1) \cdots f(e_n)$$

$$= f(x_1 - \theta) \cdots f(x_n - \theta)$$

$$\frac{d \log L(\theta)}{d\theta} = 0$$

$$\hat{\theta} = \bar{x}$$

$$\sum_{i=1}^n \frac{f'(x_i - \theta)}{f(x_i - \theta)} = 0$$

- 最大似然估计到正态分布

$$\sum_{i=1}^n \frac{f'(x_i - \theta)}{f(x_i - \theta)} = 0 \quad \hat{\theta} = \bar{x}$$

$$f'(x) = e^{x^2} \cdot (x^2)' = 2xe^{x^2}$$

$$f'(x) = 2xe^{x^2}$$

$$g(x) = \frac{f'(x)}{f(x)} \rightarrow \sum_{i=1}^n g(x_i - \theta) = 0 \rightarrow \sum_{i=1}^n g(x_i - \bar{x}) = 0$$

标准正态分布：  
 $\varphi(x) = \frac{e^{-\frac{1}{2}x^2}}{\sqrt{2\pi}}$

$$n = 2$$

$$g(x_1 - \bar{x}) + g(x_2 - \bar{x}) = 0$$

$$x_1 - \bar{x} = -(x_2 - \bar{x})$$

$$g(-x) = -g(x)$$

$$n = m + 1 \quad x_1 = \dots = x_m = -x, x_{m+1} = mx$$

$$\sum_{i=1}^n g(x_i - \bar{x}) = mg(-x) + g(mx)$$

$$g(mx) = mg(x)$$

$$g(x) = cx$$

$$f(x) = Me^{cx^2}$$

- 改变充分统计量到指数族分布

Distribution	Parameter(s) $\theta$	Natural parameter(s) $\eta$	Base measure $h(x)$	Sufficient statistic $T(x)$
Bernoulli distribution	$p$	$\ln \frac{p}{1-p}$ • This is the logit function.	1	$x$
Poisson distribution	$\lambda$	$\ln \lambda$	$\frac{1}{x!}$	$x$
exponential distribution	$\lambda$	$-\lambda$	1	$x$
normal distribution	$\mu, \sigma^2$	$\begin{bmatrix} \frac{\mu}{\sigma^2} \\ -\frac{1}{2\sigma^2} \end{bmatrix}$	$\frac{1}{\sqrt{2\pi}}$	$\begin{bmatrix} x \\ x^2 \end{bmatrix}$
gamma distribution	$\alpha, \beta$	$\begin{bmatrix} \alpha - 1 \\ -\beta \end{bmatrix}$	1	$\begin{bmatrix} \ln x \\ x \end{bmatrix}$
	$k, \theta$	$\begin{bmatrix} k - 1 \\ -\frac{1}{\theta} \end{bmatrix}$		

$$f_X(x | \theta) = h(x)g(\theta) \exp(\eta(\theta) \cdot T(x))$$

## 2.4 贝叶斯派

- 最大相互熵的求解：满足分布积分为 1 和一阶均值矩估计
  1. 初始观察分布  $m(x)$ : 这个可以是随意的观察情况，不一定必须是一个分布函数。
  2.  $\sum_{x \in S} f(x) = 1, f(x) \geq 0$ : 满足分布函数的情况，这里仅仅考虑离散的情况。
  3.  $\sum_{x \in S} t_j(x)f(x) = \mathbb{E}\{t_j(x)\} = \mu_j$ , 其中  $j \in J$ :  $t_j(x)$  是在结合上的一个测量函数，并且我们知道测量值的期望是  $\mu_j$ 。

目标:  $\operatorname{argmax}_{f(x)} RE(f(x) \| m(x))$ : 我们希望找一个函数满足上述的限制条件，并且尽可能与初始观察分布的相对熵最大。

$$RE(f(x) \| m(x)) = - \sum_{x \in S} f(x) \log \frac{f(x)}{m(x)}$$

- 应用拉格朗日乘数法:

$$L(f) = - \sum_{x \in S} f(x) \log \frac{f(x)}{m(x)} + \lambda \left( \sum_{x \in S} f(x) - 1 \right) + \sum_{j \in J} \theta_j \left( \sum_{x \in S} t_j(x) f(x) - \mu_j \right)$$

令  $L(f)$  对  $f$  的导数为 0

$$\begin{aligned} 0 &= \frac{\partial L(f)}{\partial f} = -\left(\log \frac{f(x)}{m(x)} + 1\right) + \lambda + \sum_{j \in J} \theta_j t_j(x) \\ &= -\log f(x) + \log m(x) - 1 + \lambda + \sum_{j \in J} \theta_j t_j(x) \end{aligned}$$

从而得到

$$f(x) = m(x) \exp[\lambda - 1 + \sum_{j \in J} \theta_j t_j(x)] \quad (1)$$

- 接下来应用概率求和为 1 的限制条件：

$$\begin{aligned} 1 &= \sum_{x \in S} f(x) = \sum_{x \in S} m(x) \exp\{\lambda - 1 + \sum_{j \in J} \theta_j t_j(x)\} \\ &= e^{\lambda-1} \sum_{x \in S} m(x) \exp\{\sum_{j \in J} \theta_j t_j(x)\} \end{aligned}$$

于是有

$$1 - \lambda = \log \left( \sum_{x \in S} m(x) \exp[\sum_{j \in J} \theta_j t_j(x)] \right)$$

- 代入优化表达式化简：

我们将上式的右边定义为  $b(\vec{\theta})$ , 其中  $\vec{\theta} = (\theta_1, \dots, \theta_{|J|})$ ,  $\vec{t}(x) = (t_1(x), \dots, t_{|J|}(x))$ :

$$\begin{aligned} b(\vec{\theta}) &= \log \left( \sum_{x \in S} m(x) \exp \left[ \sum_{j \in J} \theta_j t_j(x) \right] \right) \\ &= \log \left( \mathbb{E}(e^{\vec{t}(x)^T \vec{\theta}}) \right) \\ &= 1 - \lambda \end{aligned}$$

在(1)式中替换  $1 - \lambda$  得到  $f(x)$  的如下形式：

$$\begin{aligned} f(x) &= m(x) \exp[-b(\vec{\theta}) + \sum_{j \in J} \theta_j t_j(x)] \\ &= m(x) \exp[\vec{t}(x)^T \vec{\theta} - b(\vec{\theta})] \end{aligned}$$

- 最大相互熵推导出概率分布
1. 指数族概率分布
    - 这(2)就是指数族概率分布形式
  2. 正态分布
    - 只用了 1 阶矩估计：均值限制
    - 如果用了 2 阶矩估计：方差限制，就得到我们要的正态分布了！
  3. 频率派，费希尔派，到贝叶斯派，要求的事实观测越来越少，要求的经验也越来越少。

- 以少测多，见微知著

- 从频率出发生成概率基本概念
- 从经验出发总结概率基本工具
- 从推理出发奠定概率基本原理



## 2.5 K 分类最大相互熵推导 Softmax

- 限制下求分布设  $\sigma(\mathbf{x})_v$  表示  $\mathbf{x}$  属于第  $v$  个分类的概率，共有  $K$  个分类。
  - 概率分布  $\sigma(\mathbf{x})_v$  的限制条件：

$$\begin{aligned}\sigma(\mathbf{x})_v &\geq 0 \\ \sum_{v=1}^K \sigma(\mathbf{x})_v &= 1\end{aligned}$$

- 分布一致限制：分布  $\sigma(\mathbf{x})_v$  满足训练集 (Training Set) 中数据的要求，即  $\sigma(\mathbf{x})_v$  与训练集的数据分布一致：

$$\sum_{i=1}^n \sigma(\mathbf{x}(i))_v \mathbf{x}(i)_j = \sum_{i=1}^n \mathbb{I}_v(y(i)) \mathbf{x}(i)_j$$

$\mathbb{I}_u(y(i))$  为指示函数，当  $y(i) = u$  时值为 1，当  $y(i) \neq u$  时值为 0。

- 最大熵目标在每一个分类  $u$  里，任意一个特征  $j$  在属于该分类的训练数据  $\mathbf{x}(i)$  上的求和，等于所训练的模型分配给特征  $j$  的概率质量 (Probability Mass) 之和 ( $\mathbf{x}(i)$  属

于分类  $u$  的概率，在全部训练数据上求和)。这表明  $\sigma(\mathbf{x}(i))_v$  是对训练集的指标函数  $\mathbb{I}_u(y(i))$  的一个很好的近似。

$\sigma(\mathbf{x}(i))_v$  的熵定义为：

$$-\sum_{v=1}^K \sum_{i=1}^n \sigma(\mathbf{x}(i))_v \log(\sigma(\mathbf{x}(i))_v)$$

根据 1, 2, 3 条件求最优化问题，可以使用拉格朗日乘数法求解：

$$\begin{aligned} L = & \sum_{j=1}^m \sum_{v=1}^K \lambda_{v,j} \left( \sum_{i=1}^n \sigma(\mathbf{x}(i))_v \mathbf{x}(i)_j - \sum_{i=1}^n \mathbb{I}_v(y(i)) \mathbf{x}(i)_j \right) \\ & + \sum_{i=1}^n \beta_i \left( \sum_{v=1}^K \sigma(\mathbf{x})_v - 1 \right) \\ & - \sum_{v=1}^K \sum_{i=1}^n \sigma(\mathbf{x}(i))_v \log(\sigma(\mathbf{x}(i))_v) \end{aligned}$$

式中的  $L$  对  $\sigma(\mathbf{x}(i))_v$  求偏导得到:

$$\frac{\partial L}{\partial \sigma(\mathbf{x}(i))_v} = \lambda_v \mathbf{x}(i) + \beta_i - \log(\sigma(\mathbf{x}(i))_v) - 1$$

令上式等于 0:

$$\lambda_v \mathbf{x}(i) + \beta_i - \log(\sigma(\mathbf{x}(i))_v) - 1 = 0$$

解方程得到:

$$\sigma(\mathbf{x}(i))_v = e^{\lambda_v \mathbf{x}(i) + \beta_i - 1}$$

根据限制条件, 概率  $\sigma(\mathbf{x}(i))_v$  求和等于 1:

$$\sum_{v=1}^k e^{\lambda_v \mathbf{x}(i) + \beta_i - 1} = 1$$

于是得到:

$$e^{\beta_i - 1} = \frac{1}{\sum_{v=1}^k e^{\lambda_v \mathbf{x}(i) - 1}}$$

把上式中的  $e^{\beta_i - 1}$  代入:

$$\sigma(\mathbf{x}(i)) = \frac{e^{\lambda_u \mathbf{x}(i)}}{\sum_{v=1}^k e^{\lambda_v \mathbf{x}(i)}}$$

即

$$\sigma(\mathbf{x}) = \frac{e^{\lambda_u \mathbf{x}}}{\sum_{v=1}^k e^{\lambda_v \mathbf{x}}}$$

再对 Softmax 函数的每个分量做线性扩展，那么我们就得到 Log-Linear 模型

$$p(y|x; \vec{w}) = \frac{\exp\left(c + \sum_{j=1}^{|\vec{w}|} w_j f_j(x, y)\right)}{Z(x, \vec{w})},$$
$$Z(x, \vec{w}) = \sum_{y' \in Y} \exp\left(c + \sum_{j=1}^{|\vec{w}|} w_j f_j(x, y')\right)$$

如果写成向量的形式就变成了：

$$p(y|x; \vec{w}) = \frac{\exp\left(\vec{w}^T \vec{f}(x, y)\right)}{Z(x, \vec{w})}$$

对比一下 Softmax 函数，我们把 Softmax 的每个输入自变量变成了线性：

$$v(x, y) = \vec{w}^T \vec{f}(x, y) \Rightarrow$$
$$p(y|x; \vec{w}) = \sigma(v(x, y))_{y \in Y} = \frac{e^{v(x, y)}}{\sum_{y' \in Y} e^{v(x, y')}}$$

### 3 词嵌入模型

#### 3.1 词-文上下文 Word-Document

- 向量空间模型：文档模型

$$w_{x,y} = \text{tf}_{x,y} \times \log\left(\frac{N}{df_x}\right)$$

## TF-IDF

Term  $x$  within document  $y$

$\text{tf}_{x,y}$  = frequency of  $x$  in  $y$

$df_x$  = number of documents containing  $x$

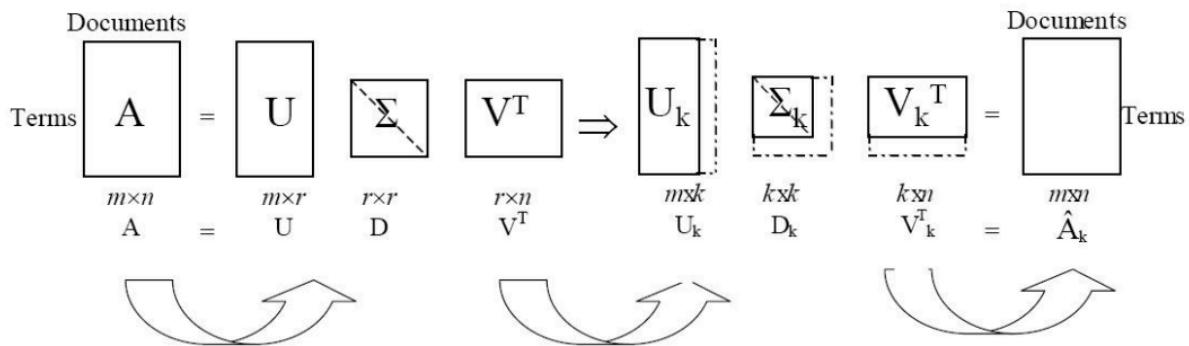
$N$  = total number of documents

	Document 1	Document 2	Document 3	Document 4	Document 5	Document 6	Document 7	Document 8
Term(s) 1	10	0	1	0	0	0	0	2
Term(s) 2	0	2	0	0	0	18	0	2
Term(s) 3	0	0	0	0	0	0	0	2
Term(s) 4	6	0	0	4	6	0	0	0
Term(s) 5	0	0	0	0	0	0	0	2
Term(s) 6	0	0	1	0	0	1	0	0
Term(s) 7	0	1	8	0	0	0	0	0
Term(s) 8	0	0	0	0	0	3	0	0

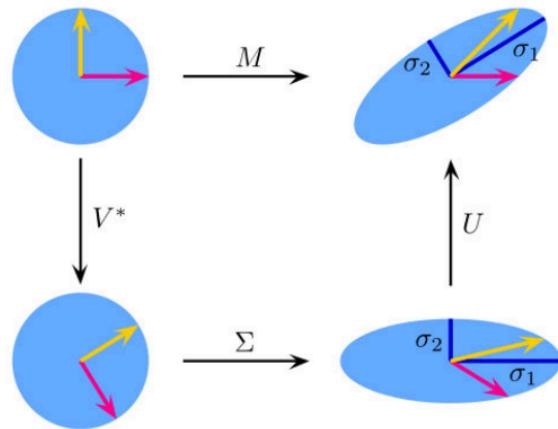
↑  
Document Vector

Word Vector  
(Passage Vector)

- 向量空间矩阵的 SVD



- Singular Value Decomposition (SVD)



$$M = U \cdot \Sigma \cdot V^*$$

- 隐性语义索引 Latent Semantic Indexing (LSI)

**A**

$$\mathbf{M} = \begin{matrix} & D_1 & D_2 & D_3 & D_4 & D_5 & D_6 & \cdots & D_n \\ T_1 & 0.00060 & 0.00012 & 0.00003 & 0.00003 & 0.0033 & 0.00048 & \cdots & a_{1n} \\ T_2 & 0 & 0 & 0 & 0 & 0 & 0 & \cdots & a_{2n} \\ T_3 & 0 & 2.98862 & 0 & 0 & 0 & 1.49431 & \cdots & a_{3n} \\ T_4 & 0 & 0 & 0 & 13.32555 & 0 & 0 & \cdots & a_{4n} \\ T_5 & 0 & 0 & 0 & 0 & 0 & 0 & \cdots & a_{5n} \\ T_6 & 1.03442 & 1.03442 & 0 & 0 & 0 & 3.10326 & \cdots & a_{6n} \\ \vdots & \ddots & \vdots \\ T_m & a_{m1} & a_{m2} & a_{m3} & a_{m4} & a_{m5} & a_{m6} & \cdots & a_{mn} \end{matrix}$$

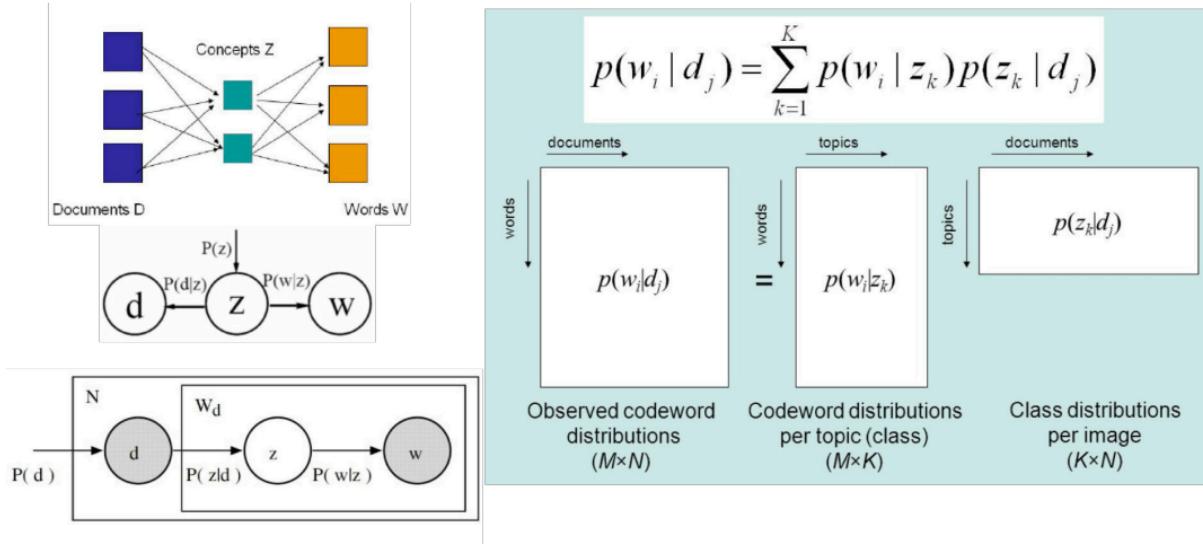
**B**

$$\mathbf{U}_k = \boxed{\begin{matrix} C_1 & C_2 & C_3 & \cdots & C_m \\ T_1 & a_{11} & a_{12} & a_{13} & \cdots & a_{1m} \\ T_2 & a_{21} & a_{22} & a_{23} & \cdots & a_{2m} \\ T_3 & a_{31} & a_{32} & a_{33} & \cdots & a_{3m} \\ T_4 & a_{41} & a_{42} & a_{43} & \cdots & a_{4m} \\ T_5 & a_{51} & a_{52} & a_{53} & \cdots & a_{5m} \\ T_6 & a_{61} & a_{62} & a_{63} & \cdots & a_{6m} \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ T_m & a_{m1} & a_{m2} & a_{m3} & \cdots & a_{mm} \end{matrix}}$$

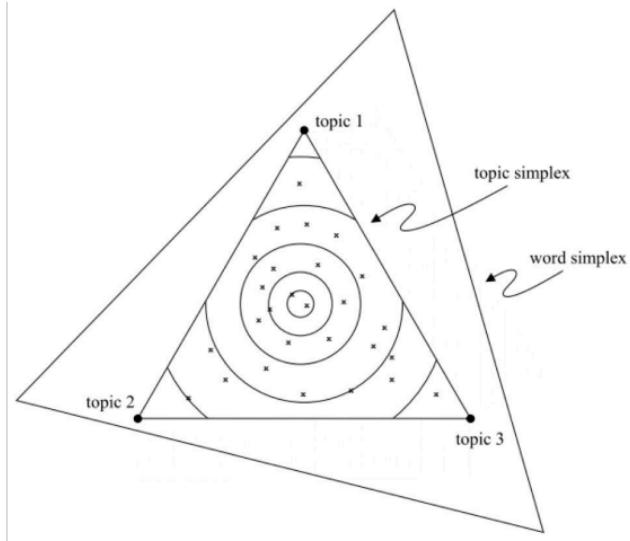
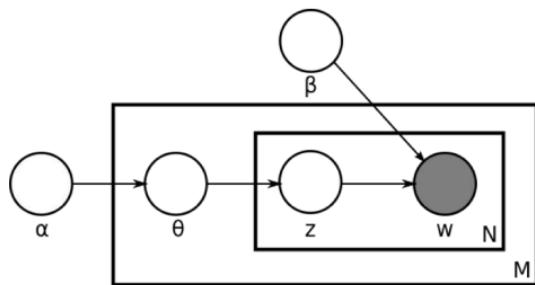
$$\Sigma_k = \boxed{\begin{matrix} D_1 & D_2 & D_3 & \cdots & D_n \\ T_1 & a_{11} & 0 & 0 & \cdots & 0 \\ T_2 & 0 & a_{22} & 0 & \cdots & 0 \\ T_3 & 0 & 0 & a_{33} & \cdots & 0 \\ T_4 & 0 & 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ T_m & 0 & 0 & 0 & \cdots & a_{nn} \end{matrix}}$$

$$\mathbf{V}^T = \boxed{\begin{matrix} D_1 & D_2 & D_3 & \cdots & D_n \\ C_1 & a_{11} & a_{12} & a_{13} & \cdots & a_{1n} \\ C_2 & a_{21} & a_{22} & a_{23} & \cdots & a_{2n} \\ C_3 & a_{31} & a_{32} & a_{33} & \cdots & a_{3n} \\ C_4 & a_{41} & a_{42} & a_{43} & \cdots & a_{4n} \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ C_n & a_{n1} & a_{n2} & a_{n3} & \cdots & a_{nn} \end{matrix}}$$

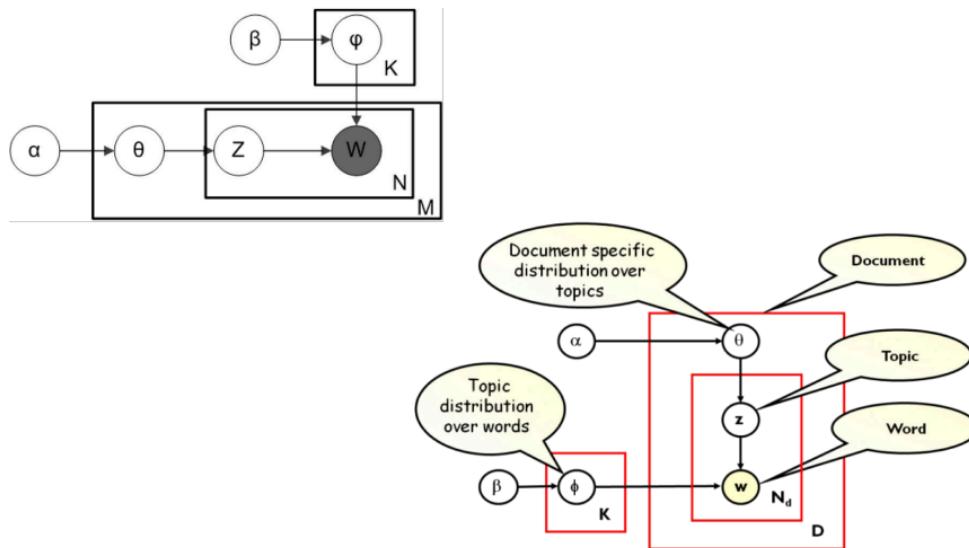
- pLSA: probabilistic Latent Semantic Analytics



- LDA: Latent Dirichlet allocation

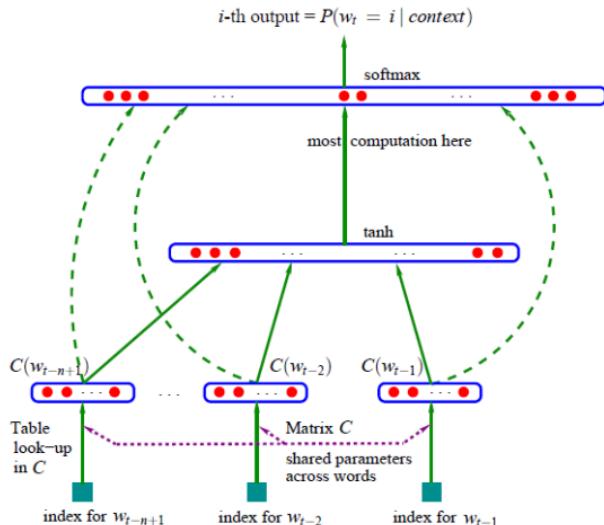


- LDA: Latent Dirichlet allocation

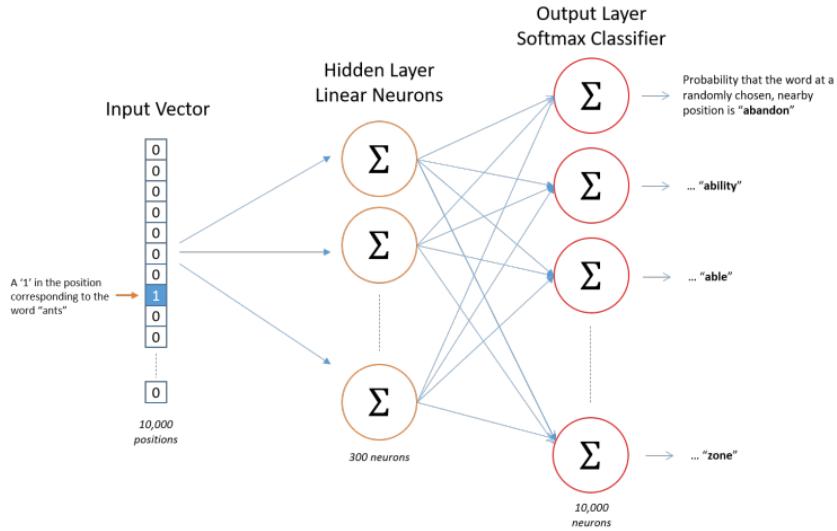


### 3.2 词上下文 Word-Neighboring Word

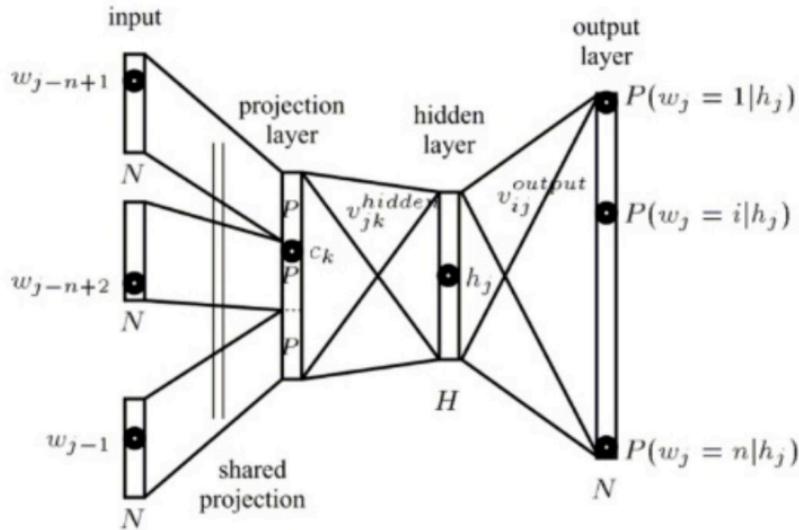
- Neural Net Language Model 2003



- Softmax Network

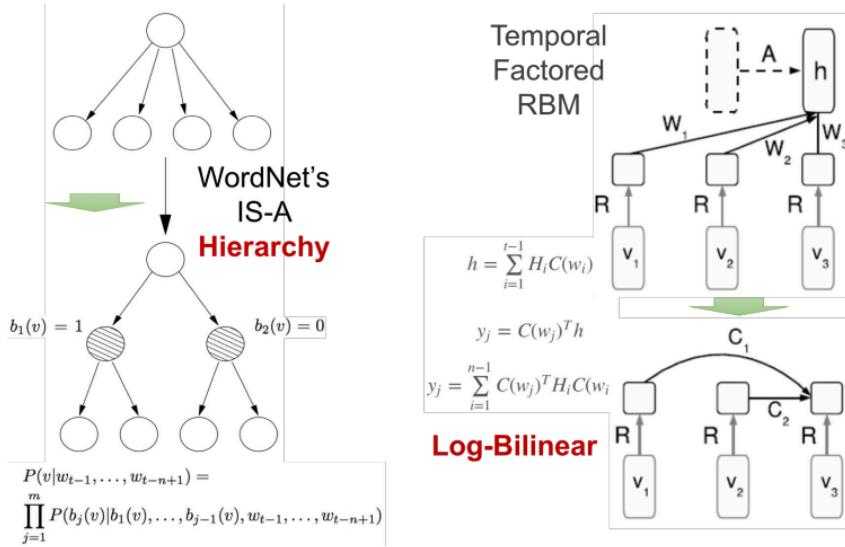


- NNLM 2006: 3 层模型

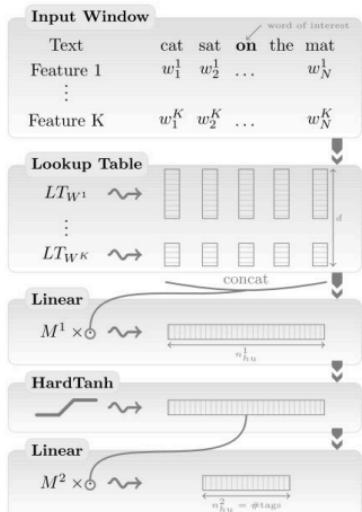


- NNLM 小结
1. 开山之作：2001 年发在 NIPS，2003 同名发到 JMLR
  2. 三层的神经网络来构建语言模型
  3. 一词多义有待解决，9 年后 Huang 模型搞定
  4. 降低参数个数，比如用 RNN, Mikolov 搞定
  5. 提出计算量问题，开启 Softmax 近似计算

- M&H 的 Hierarchical Log-BiLinear (HLBL) 模型 2008

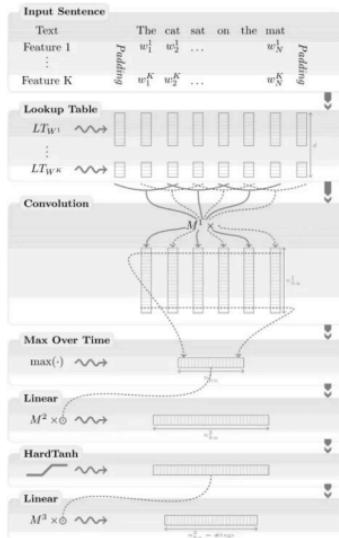


- M&H HLBL 小结
  1. Softmax 层变成平衡二叉树输出，训练速度提升
  2. Factored RBM(Temporal Factored RBM) 简化到 Log-BiLinear
    - 相邻词向量矩阵对应的词向量表示，通过链接矩阵和当前词关联，而关联关系使用能量函数表示
  3. Log-BiLinear 合并隐藏层和输出成，简化网络，训练速度提升 (Contrastive Divergence 算法)
  4. Hierarchical Log-BiLinear 限制对上下文 Context 的关系矩阵为对角矩阵
  5. 启发式，随机初始化二叉树 (Bagging 思想)
- Collobert and Weston (SENNER) 模型 2008: Window 窗口



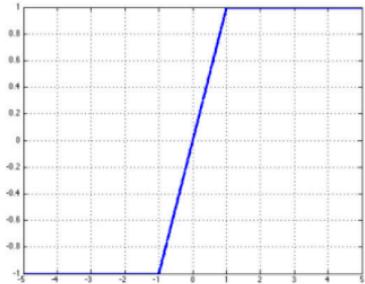
HardTanh 取代了 Tanh

- Collobert and Weston (SENNNA) 模型 2008: Sentence 句子



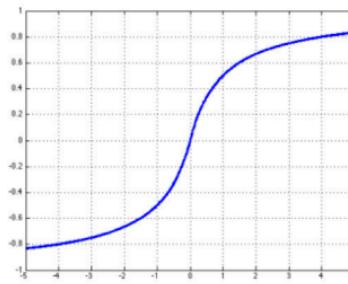
CNN 压缩句子

- HardTanh



$$\text{hardtanh}(z) = \begin{cases} -1 & : z < -1 \\ z & : -1 \leq z \leq 1 \\ 1 & : z > 1 \end{cases}$$

$$\text{hardtanh}'(z) = \begin{cases} 1 & : -1 \leq z \leq 1 \\ 0 & : \text{otherwise} \end{cases}$$



$$\text{softsign}(z) = \frac{z}{1 + |z|}$$

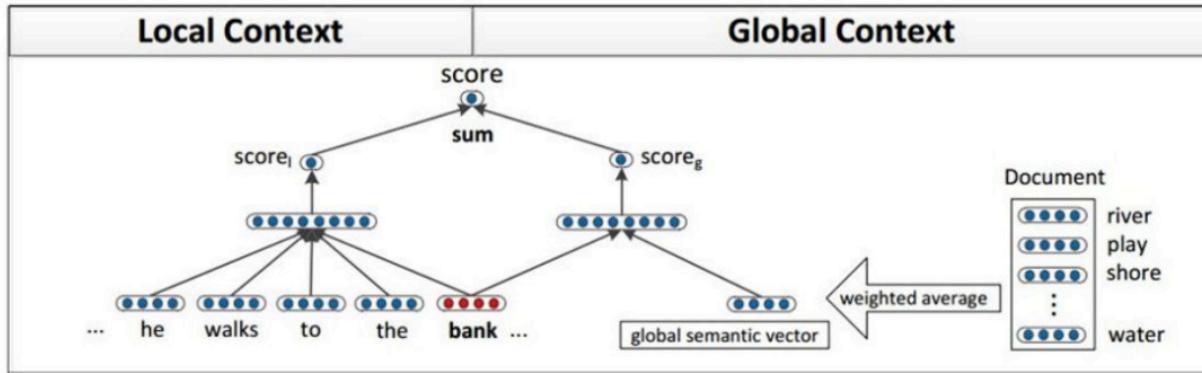
$$\text{softsign}'(z) = \frac{\text{sgn}(z)}{(1+z)^2}$$

- C&W SENNA 小结
  1. 专注到 NLP 应用：词性标注、命名实体识别、短语识别、语义角色标注
  2. 近似计算：使用 pair-wise 的方法训练词向量

$$\sum_{x \in \mathfrak{X}} \sum_{w \in \mathfrak{D}} \max\{0, 1 - f(x) + f(x^{(w)})\}$$

3. 输出层只有一个节点，而不是 NNLM 的  $|V|$  个节点，降低计算复杂度
4. 大写开头的单词和小写单词当作同一个词处理
5. 维基百科英文语料和路透社语料中一共训练了 7 周

- Eric Huang's Model

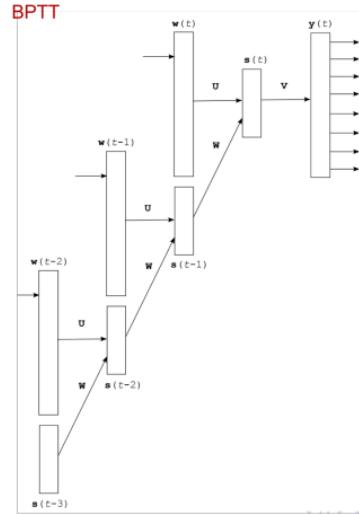
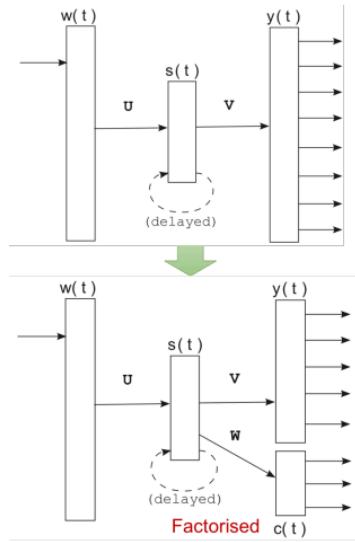


- Clustering word base on context and retrain the model

- Eric Huang's Model 小结

1. 基于 C&W 的基础的改进 (pair-wise 目标函数)
2. 使用全文信息辅助已有的局部信息
  - 全局信息：文章中所有词的词向量求个加权平均 (权重是词的 idf)
3. 使用多个词向量来表示多义词
  - 词的上下文向量：各 5 个词 (共 10 个词) 的词向量做加权平均
  - 对上下文向量做 k-means 聚类，根据聚类结果给每个词打上标签
  - 重新训练词向量

- RNN-LM: Mikolov 2010



- RNN-LM 模型

$$x(t) = w(t) + s(t-1) \quad (1)$$

$$s_j(t) = f \left( \sum_i x_i(t) u_{ji} \right) \quad (2)$$

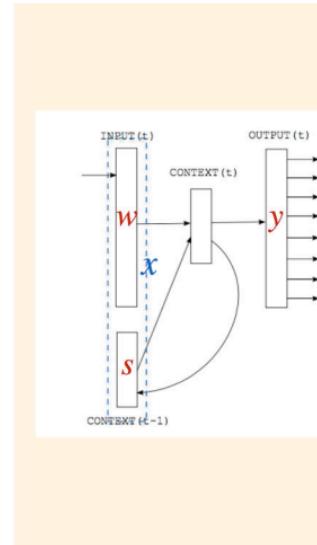
$$y_k(t) = g \left( \sum_j s_j(t) v_{kj} \right) \quad (3)$$

where  $f(z)$  is sigmoid activation function:

$$f(z) = \frac{1}{1 + e^{-z}} \quad (4)$$

and  $g(z)$  is softmax function:

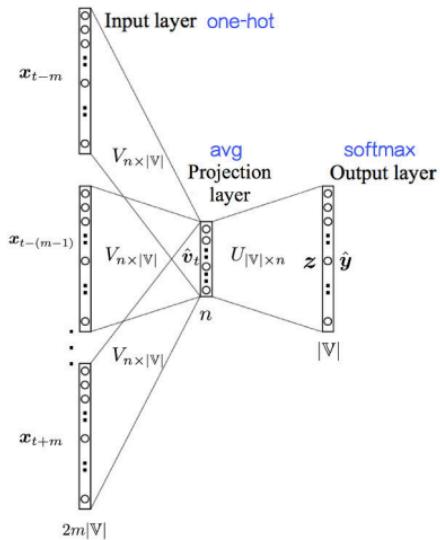
$$g(z_m) = \frac{e^{z_m}}{\sum_k e^{z_k}} \quad (5)$$



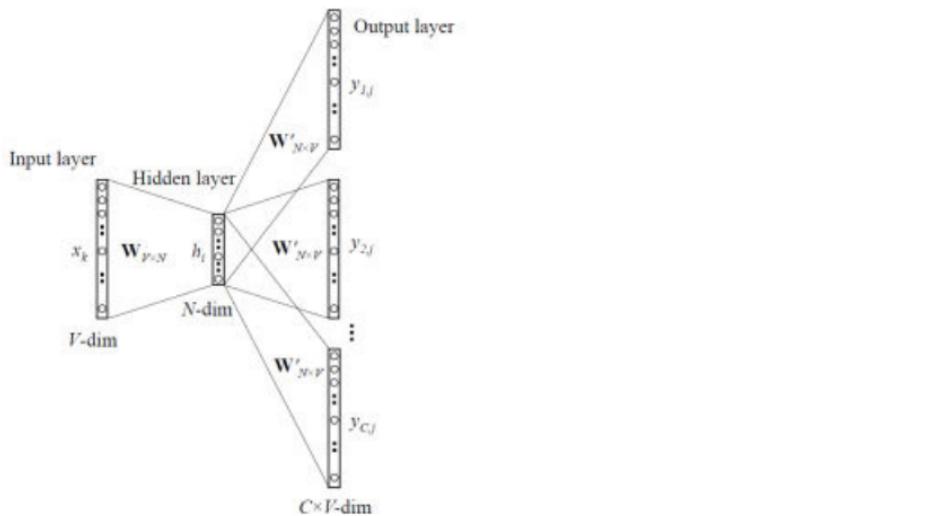
- RNN-LM 小结

1. RNN 可以解决 Long-Term Dependency 问题
2. RNN 是深度网络可以发现更多词与词的 Pattern
3. RNN 让训练变得复杂，训练困难
4. RNN-LM 采用分组计算 Softmax
5. RNN-LM 隐藏层采用了 Sigmoid，均值不在 0 上下

- CBOW (Continuous Bag of words): 给定周围词，预测中心词



- Skip-Gram: 给定词，预测周围词 (随机选一个周围词)



- CBOW

$$J_\theta = \frac{1}{T} \sum_{t=1}^T \log p(w_t | w_{t-n}, \dots, w_{t-1}, w_{t+1}, \dots, w_{t+n})$$

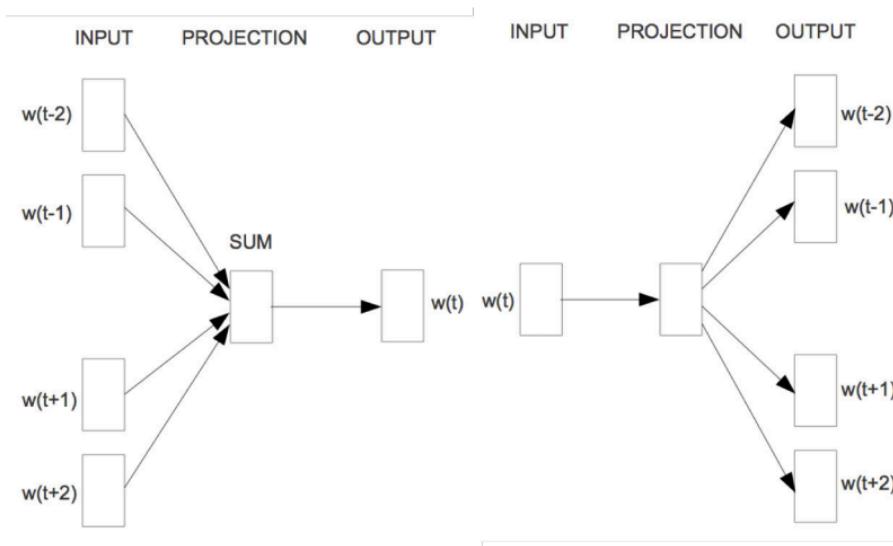
$$p(w_t | w_{t-1}, \dots, w_{t-n+1}) = \frac{\exp(h^\top v'_{w_t})}{\sum_{w_i \in V} \exp(h^\top v'_{w_i})}$$

- Skip-Gram

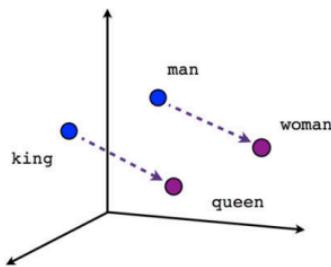
$$J_\theta = \frac{1}{T} \sum_{t=1}^T \sum_{-n \leq j \leq n, j \neq 0} \log p(w_{t+j} | w_t)$$

$$p(w_{t+j} | w_t) = \frac{\exp(h^\top v'_{w_{t+j}})}{\sum_{w_i \in V} \exp(h^\top v'_{w_i})}$$

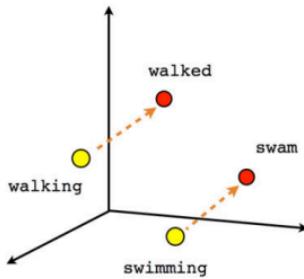
- Word2Vec



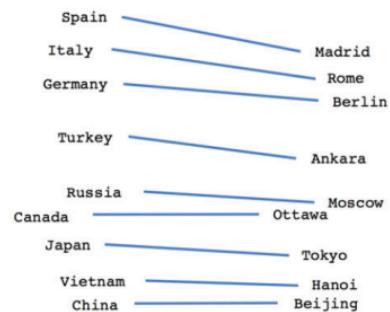
- Word2Vec 效果



Male-Female



Verb tense

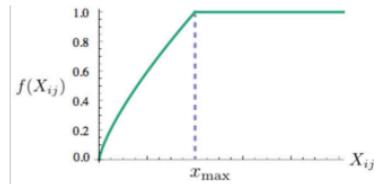


Country-Capital

- Word2Vec 小结
  1. CBOW 和 Skip-Gram 基本就是 Softmax 网络
  2. 目标是负的 Log 损失 (交互熵)
  3. 中间是线性单元，而不是 Tanh
  4. 输入是上下文（周围词），而不是上文
  5. 计算速度快，容易重现
  6. 近似计算 Softmax
    - Hierarchical Softmax （同 Hierarchical NNLM）
    - Noise Contrastive Estimation
    - Negative Sampling (简化自 NCE)

- Glove (Global Vectors for word representation): Jeffrey 和 Socher

Probability and Ratio	$k = solid$	$k = gas$	$k = water$	$k = fashion$
$P(k ice)$	$1.9 \times 10^{-4}$	$6.6 \times 10^{-5}$	$3.0 \times 10^{-3}$	$1.7 \times 10^{-5}$
$P(k steam)$	$2.2 \times 10^{-5}$	$7.8 \times 10^{-4}$	$2.2 \times 10^{-3}$	$1.8 \times 10^{-5}$
$P(k ice)/P(k steam)$	8.9	$8.5 \times 10^{-2}$	1.36	0.96

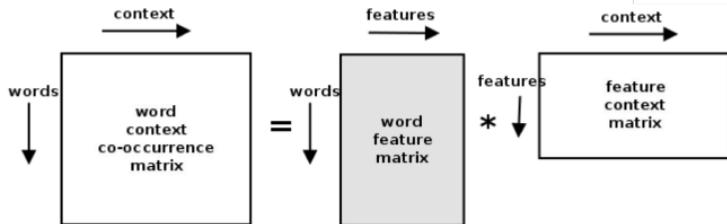


共现权重

$$f(x) = \begin{cases} (x/x_{\max})^{\alpha} & \text{if } x < x_{\max} \\ 1 & \text{otherwise} \end{cases}$$

$$\sum_{i,j=1}^V f(X_{i,j})(u_i^T v_j + b_i + c_j - \log X_{i,j})^2$$

加权最小二乘目标



- Glove 目标定义

加权最小二乘目标：

$$J = \sum_{i,j=1}^V f(X_{ij}) \left( w_i^T \tilde{w}_j + b_i + \tilde{b}_j - \log X_{ik} \right)^2$$

定义条件概率 ( $i, j$  的词频) :

$$P_{ij} = P(j|i) = \frac{X_{ij}}{X_i}$$

条件概率比值 (给定  $k$  后,  $i, j$  的条件概率比) :

$$F(w_i, w_j, \tilde{w}_k) = \frac{P_{ik}}{P_{jk}}$$

假设存在词向量**加性**关系：

$$F(w_i - w_j, \tilde{w}_k) = \frac{P_{ik}}{P_{jk}}$$

假设存在词向量**乘性**关系：

$$F((w_i - w_j)^T \tilde{w}_k) = \frac{P_{ik}}{P_{jk}}$$

并且**统一函数**关系：

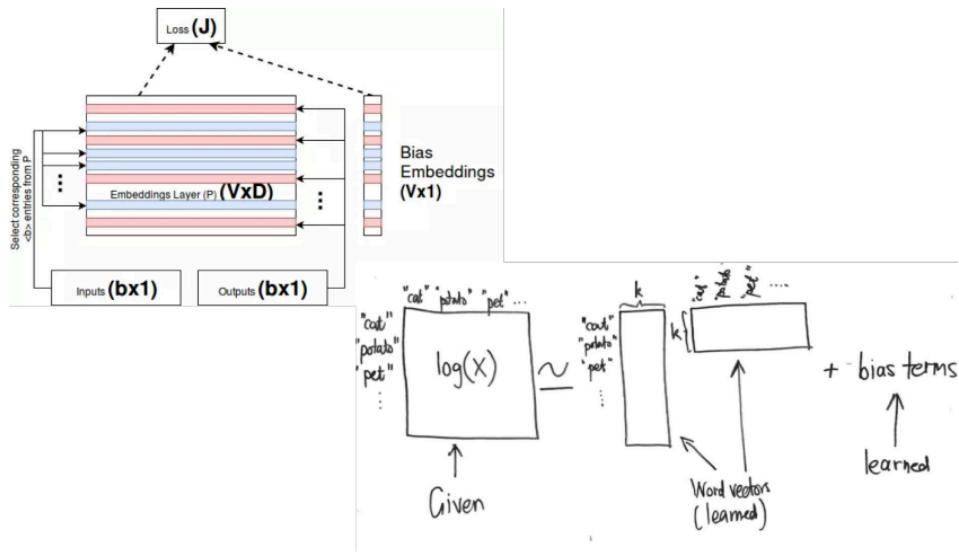
$$F((w_i - w_j)^T \tilde{w}_k) = \frac{F(w_i^T \tilde{w}_k)}{F(w_j^T \tilde{w}_k)}$$

$$F(w_i^T \tilde{w}_k) = P_{ik} = \frac{X_{ik}}{X_i}$$

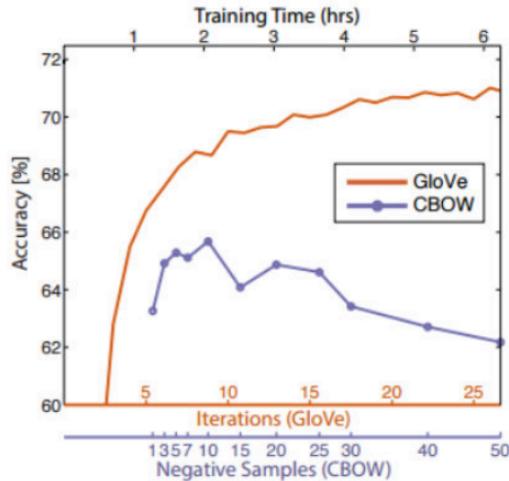
**统一函数**为指数函数：

$$w_i^T \tilde{w}_k + b_i + \tilde{b}_k = \log X_{ik}$$

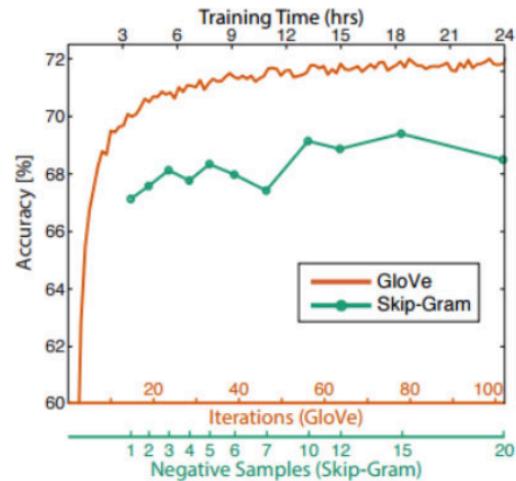
- Glove



- Glove vs Word2Vec

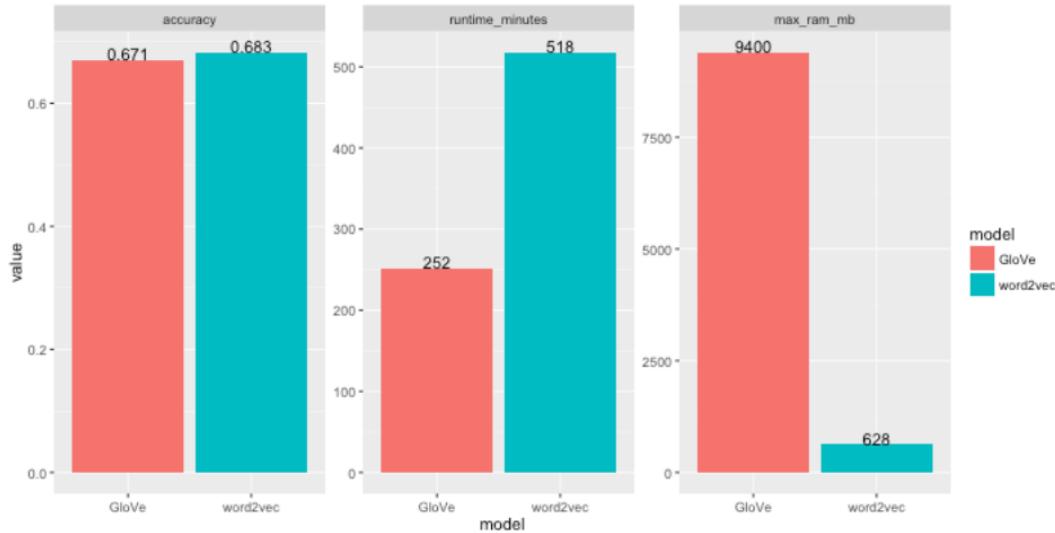


(a) GloVe vs CBOW



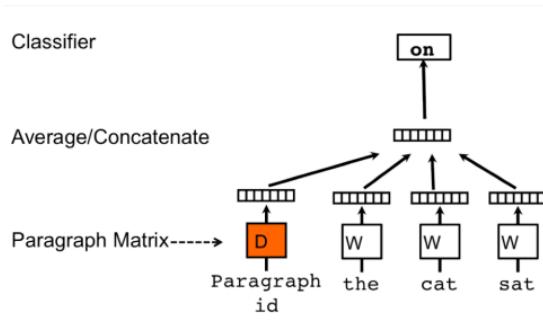
(b) GloVe vs Skip-Gram

- Glove vs Word2Vec



- Glove 小结
  1. 是 LSI 和 Log Bi-Linear 的回归的融合
  2. 吸收了 Huang's Model 全局思想, 分为 Global Matrix Factorization 和 Local Context Window
  3. 引入共现权重
  4. 训练速度比 Word2Vec 要快, 效果不比 Word2Vec 差, 就是内存占用高于 Word2Vec
  5. 在 word analogies, word similarity, NER 任务中取得不错的效果
- Paragraph2Vec

# Paragraph Vector



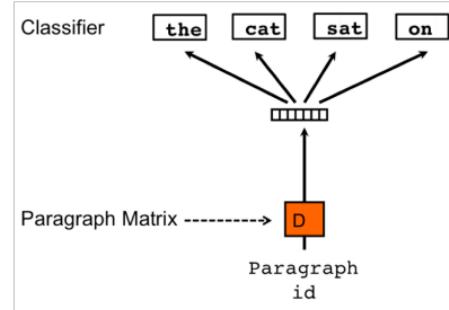
Average/Concatenate

Paragraph Matrix -----&gt;

Paragraph id

PV-DM

Distributed Memory



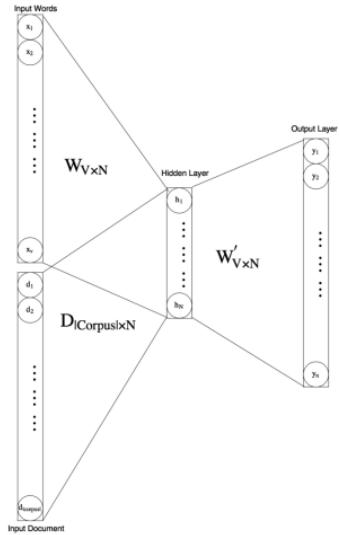
Paragraph Matrix -----&gt;

Paragraph id

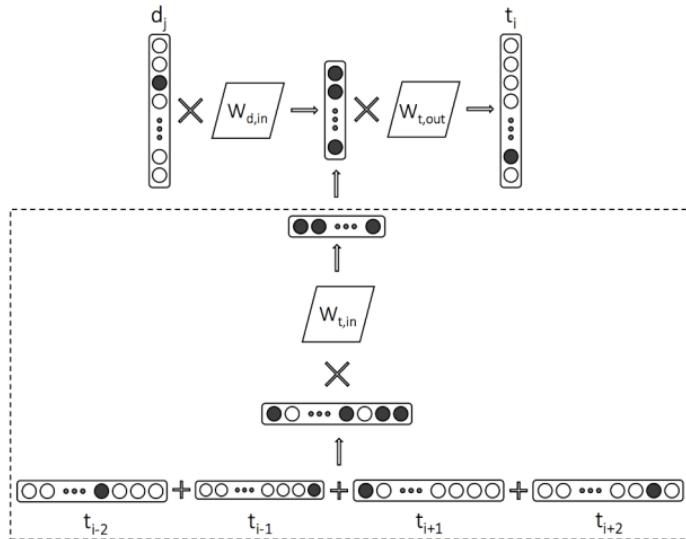
PV-DBOW

Distributed Bag of Words

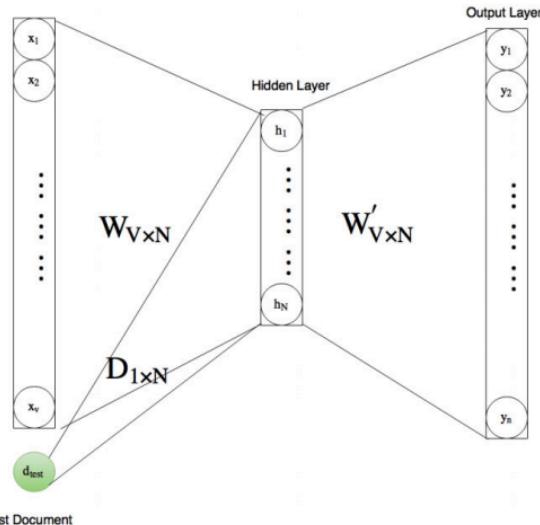
- Paragraph2Vec 训练



- Paragraph2Vec 窗口操作



- Paragraph2Vec 推理

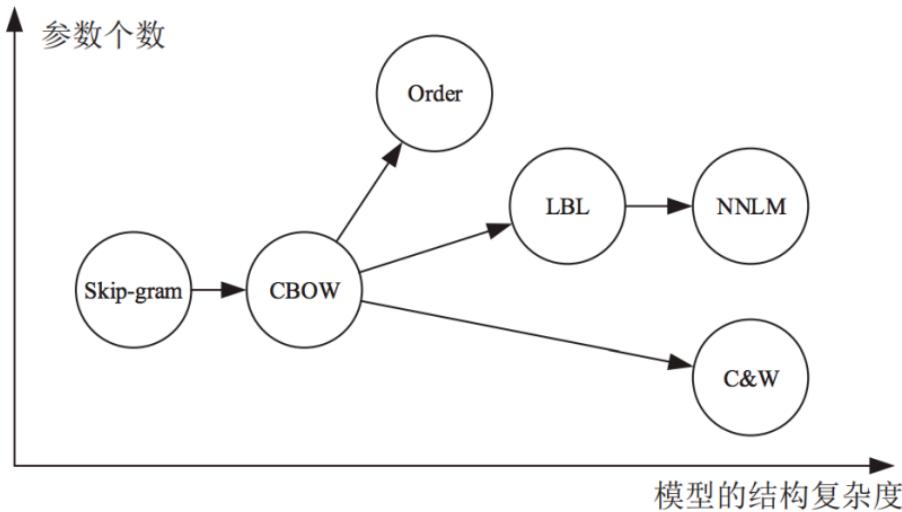


- Paragraph2Vec 对比

Model	Error rate
Vector Averaging	10.25%
Bag-of-words	8.10 %
Bag-of-bigrams	7.28 %
Weighted Bag-of-bigrams	5.67%
Paragraph Vector	<b>3.82%</b>

- Paragraph2Vec 小结
  - 1. Paragraph2Vec 保留了部分词的顺序和语义空间
    - 固定窗口输入
    - 滑动窗口到整个文本
    - Paragraph 共享一个 ID, 类似 Bagging
  - 2. 拼接处理效果优于求平均
  - 3. PV-DM 的效果要优于 PV-DBOW
  - 4. 效果略微优于 n-gram 的 BoW 模型

- 各种模型的比较



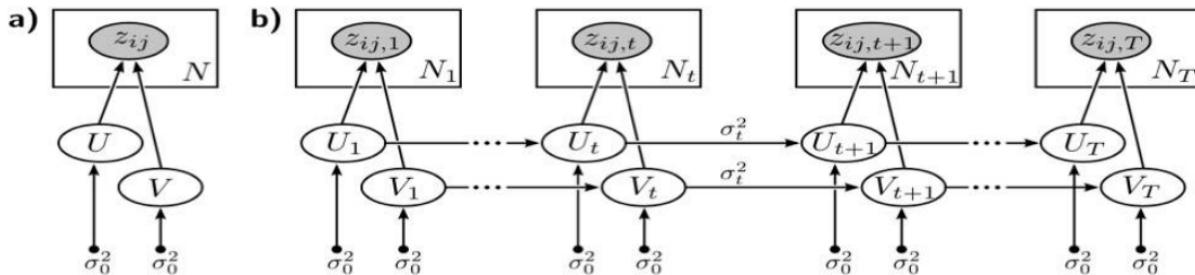
- 各种模型的比较

名称	上下文	上下文与目标词之间的建模 (技术手段)
LSA/LSI [56]	文档	
HAL [70]	词	
GloVe [90]	词	矩阵
Jones & Mewhort [45]	n-gram	
Brown Clustering [12]	词	聚类
Skip-gram [73]	词	
CBOW [73]	n-gram (加权)	
Order (2.2.7 小节)	n-gram (线性组合)	
LBL [79]	n-gram (线性组合)	神经网络
NNLM [7]	n-gram (非线性组合)	
C&W [17]	n-gram (非线性组合)	

- Word Embedding VS Topic Model

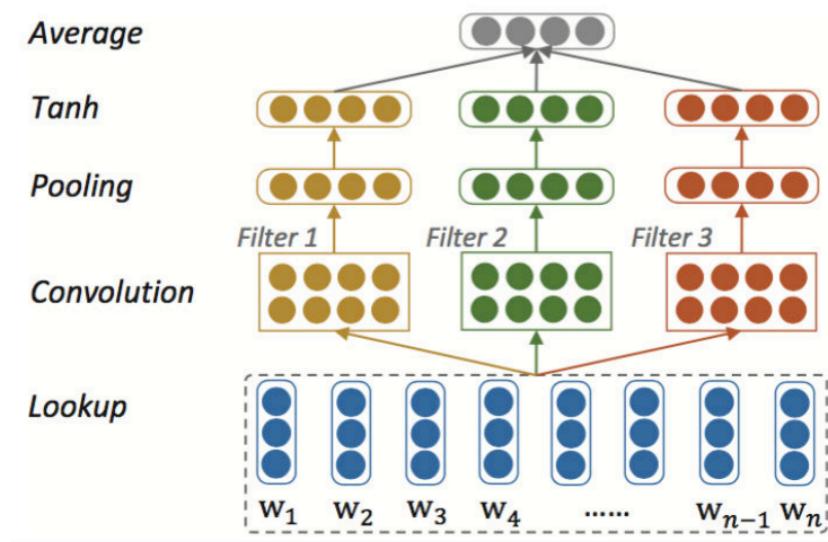
1. 词嵌入在相似度方面和词类比（word analogy）现象由于主题模型
2. 主题模型丢失了词的顺序信息
3. 主题模型的可解释性高一些，词嵌入难以解释，也没有统一的表示（例如，概率分布）
4. 词嵌入是浅层神经网络，而主题模型是浅层概率图模型
5. 词嵌入训练太慢

- 词嵌入和概率图模型的融合



Bayesian Skip-Gram

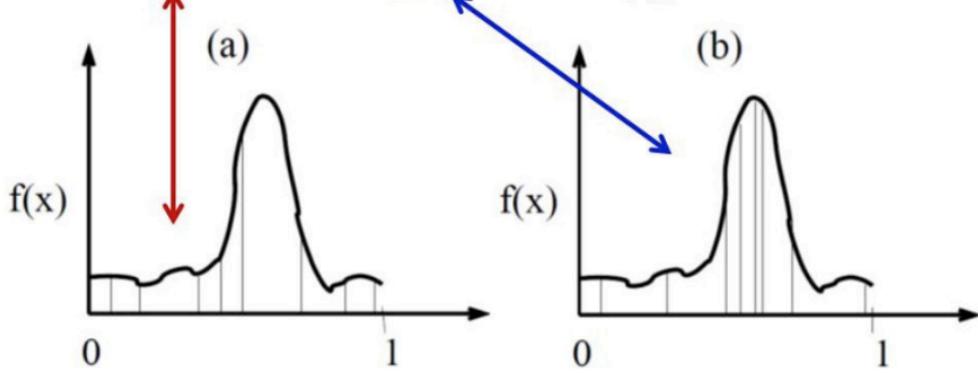
- CNN 直接词特征提取



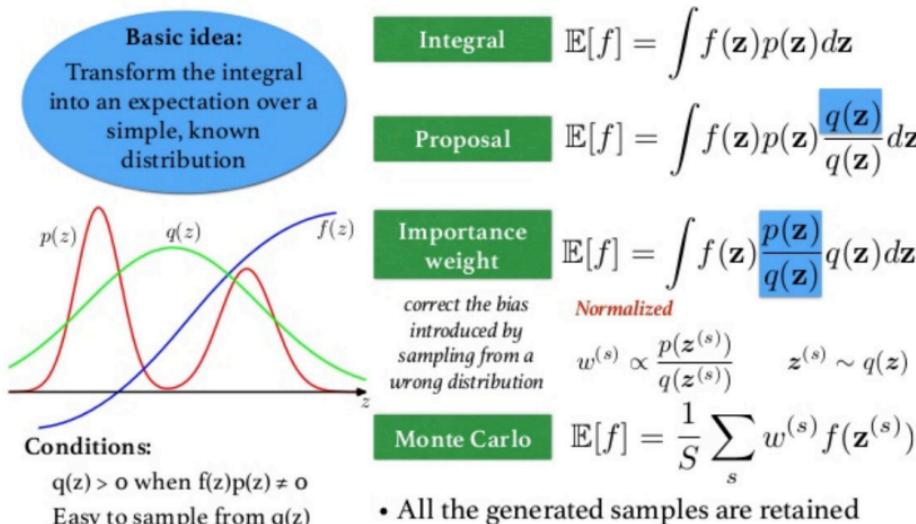
### 3.3 Softmax 近似计算

- Importance Sampling: 未知分布

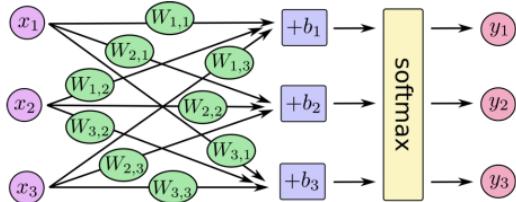
$$A = \int_a^b f(x) dx = \int_a^b \frac{f(x)}{\rho(x)} \rho(x) dx \equiv \left\langle \frac{f(x)}{\rho(x)} \right\rangle_{\rho(x)} \approx \frac{1}{M} \sum_{i=1}^M \frac{f(x_i)}{\rho(x_i)}$$



- Importance Sampling: 未知分布上的期望



- Softmax 求解看成是求期望



$$J_\theta = -\log \frac{\exp(h^\top v'_w)}{\sum_{w_i \in V} \exp(h^\top v'_{w_i})}$$

$$J_\theta = -h^\top v'_w + \log \sum_{w_i \in V} \exp(h^\top v'_{w_i})$$

$$h^\top v'_w = -\mathcal{E}(w)$$

$$J_\theta = \mathcal{E}(w) + \log \sum_{w_i \in V} \exp(-\mathcal{E}(w_i))$$

$$\nabla_\theta J_\theta = \nabla_\theta \mathcal{E}(w) + \nabla_\theta \log \sum_{w_i \in V} \exp(-\mathcal{E}(w_i))$$

$$\nabla_\theta J_\theta = \nabla_\theta \mathcal{E}(w) + \frac{1}{\sum_{w_i \in V} \exp(-\mathcal{E}(w_i))} \sum_{w_i \in V} \exp(-\mathcal{E}(w_i)) \nabla_\theta (-\mathcal{E}(w_i))$$

$$\nabla_\theta J_\theta = \nabla_\theta \mathcal{E}(w) + \sum_{w_i \in V} \frac{\exp(-\mathcal{E}(w_i))}{\sum_{w_i \in V} \exp(-\mathcal{E}(w_i))} \nabla_\theta (-\mathcal{E}(w_i))$$

$$\nabla_\theta J_\theta = \nabla_\theta \mathcal{E}(w) + \sum_{w_i \in V} P(w_i) \nabla_\theta (-\mathcal{E}(w_i))$$

$$\nabla_\theta J_\theta = \nabla_\theta \mathcal{E}(w) - \sum_{w_i \in V} P(w_i) \nabla_\theta \mathcal{E}(w_i)$$

$$\sum_{w_i \in V} P(w_i) \nabla_\theta \mathcal{E}(w_i) = \mathbb{E}_{w_i \sim P} [\nabla_\theta \mathcal{E}(w_i)]$$

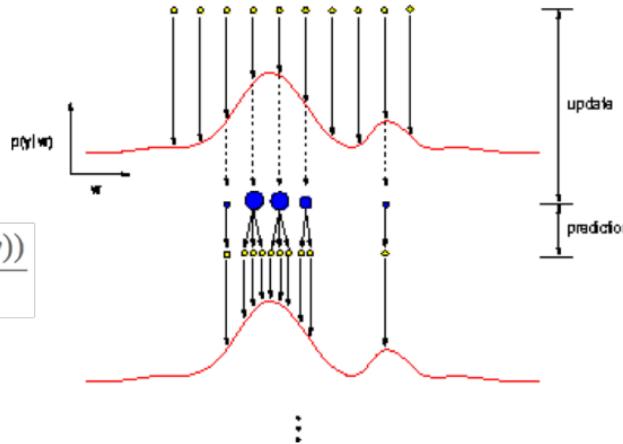
- Softmax 的 Importance Sampling 求解

$$\mathbb{E}_{w_i \sim P} [\nabla_{\theta} \mathcal{E}(w_i)] \approx \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} \mathcal{E}(w_i).$$

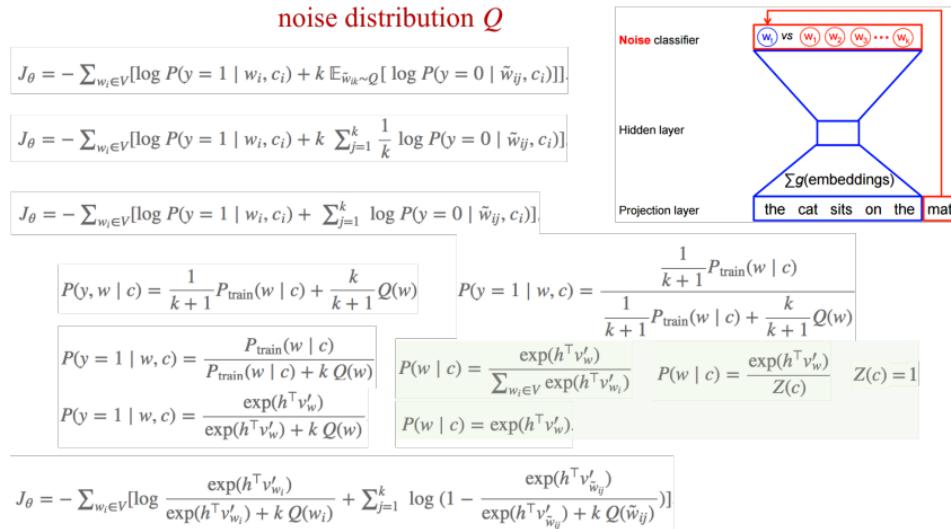
$$\mathbb{E}_{w_i \sim P} [\nabla_{\theta} \mathcal{E}(w_i)] \approx \frac{1}{R} \sum_{i=1}^m r(w_i) \nabla_{\theta} \mathcal{E}(w_i)$$

$$r(w) = \frac{\exp(-\mathcal{E}(w))}{Q(w)}$$

$$R = \sum_{j=1}^m r(w_j).$$



- Noise Contrastive Estimation : Gutmann and Hyvärinen 2010



- Negative Sampling: Mikolov 2013

$$J_\theta = - \sum_{w_i \in V} [\log \frac{\exp(h^\top v'_{w_i})}{\exp(h^\top v'_{w_i}) + k Q(w_i)} + \sum_{j=1}^k \log (1 - \frac{\exp(h^\top v'_{\tilde{w}_j})}{\exp(h^\top v'_{\tilde{w}_j}) + k Q(\tilde{w}_j)})]$$

$p(\text{professionals} | \text{xing, enables, to, grow})$

sigmoid  
negative sampling

$$P(y = 1 | w, c) = \frac{\exp(h^\top v'_w)}{\exp(h^\top v'_w) + k Q(w)}$$

$$k Q(w) = 1$$

$$k = |V|$$

$$P(y = 1 | w, c) = \frac{1}{1 + \exp(-h^\top v'_w)}$$

$$J_\theta = - \sum_{w_i \in V} [\log \frac{1}{1 + \exp(-h^\top v'_{w_i})} + \sum_{j=1}^k \log (\frac{1}{1 + \exp(h^\top v'_{\tilde{w}_j})})]$$

$$J_\theta = - \sum_{w_i \in V} [\log \sigma(h^\top v'_{w_i}) + \sum_{j=1}^k \log \sigma(-h^\top v'_{\tilde{w}_j})]$$

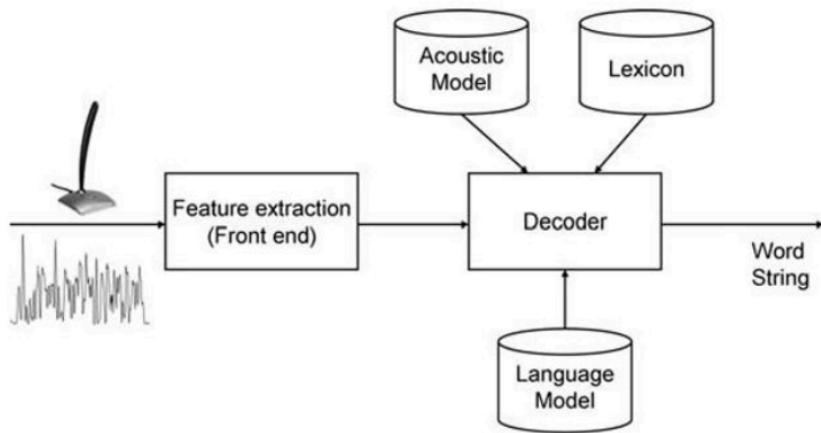
$$\sigma(x) = \frac{1}{1 + \exp(-x)}$$

- Softmax 近似小结
  1. NCE 要比 IS 稳定：NCE 利用了 2 分类结果
  2. NS 是 NCE 一种近似

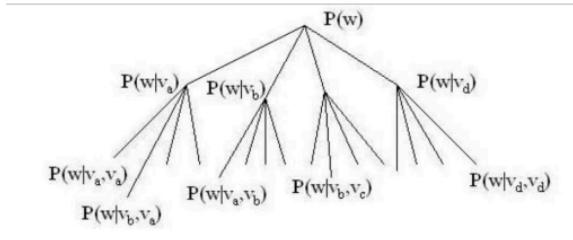
## 4 语言模型

### 4.1 从概率模型到神经网络模型

- 语言模型



- N-Gram 概率模型



## Word-level unigrams

Text	Token Sequence	Token Value
One Two Three Four	1	One
One <u>Two</u> Three Four	2	Two
One Two <u>Three</u> Four	3	Three
One Two Three <u>Four</u>	4	Four

## Word-level bigrams

Text	Token Sequence	Token Value
One Two Three Four	1	One Two
One <u>Two</u> Three Four	2	Two Three
One Two <u>Three</u> Four	3	Three Four

## Word-level trigrams

Text	Token Sequence	Token Value
One Two Three Four	1	One Two Three
One <u>Two</u> Three Four	2	Two Three Four

$$\hat{p}(w_a) = \frac{c(w_a)}{N}$$

$$\hat{p}(w_b|w_a) = \frac{c(w_a, w_b)}{\sum_{w_b} c(w_a, w_b)} \approx \frac{c(w_a, w_b)}{c(w_a)}$$

- N-Gram 概率模型

Alice was beginning to get very tired of sitting by her sister on the bank, and of having nothing to do: once or twice she had peeped into the book her sister was reading, but it had no pictures or conversations in it, 'and what is the use of a book,' thought Alice 'without pictures or conversation?'

$$P(w_{i+1} = \text{of} \mid w_i = \text{tired}) = 1$$

$$P(w_{i+1} = \text{of} \mid w_i = \text{use}) = 1$$

$$P(w_{i+1} = \text{sister} \mid w_i = \text{her}) = 1$$

$$P(w_{i+1} = \text{beginning} \mid w_i = \text{was}) = 1/2$$

$$P(w_{i+1} = \text{reading} \mid w_i = \text{was}) = 1/2$$

$$P(w_{i+1} = \text{bank} \mid w_i = \text{the}) = 1/3$$

$$P(w_{i+1} = \text{book} \mid w_i = \text{the}) = 1/3$$

$$P(w_{i+1} = \text{use} \mid w_i = \text{the}) = 1/3$$

$$P(X|Y) = \frac{P(X, Y)}{P(Y)}$$

- 语言模型的评价: 困惑度 Perplexity

熵(Entropy)的计算 :

$$H = - \lim_{m \rightarrow \infty} \frac{1}{m} \sum_{w_1, w_2, \dots, w_m} (P(w_1, w_2, \dots, w_m) \log_2 P(w_1, w_2, \dots, w_m))$$

所有情况 (全概率公式) :

$$H = - \lim_{m \rightarrow \infty} \frac{1}{m} \log_2 P(w_1, w_2, \dots, w_m)$$

有限近似 :

$$\hat{H} = - \frac{1}{m} \log_2 P(w_1, w_2, \dots, w_m)$$

困惑度(Perplexity)的计算 :

$$PP = \hat{P}(w_1, w_2, \dots, w_m)^{-\frac{1}{m}}$$

困惑度和熵的关系 :

$$PP = 2^{\hat{H}}$$

- N-Gram 模型问题

1. 概率为 0 的问题

- Laplace (Add-One) Smoothing
- Good-Turing
- Interpolation
- Backoff
- Kneser-Ney

2. 马尔可夫假设不满足问题

- N-Gram 内在问题
- Long-Term Dependency 问题

- 概率为 0 的问题

- When we have sparse statistics:

$P(w | \text{denied the})$

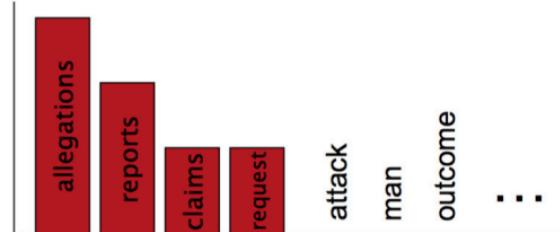
3 allegations

2 reports

1 claims

1 request

7 total



- Steal probability mass to generalize better

$P(w | \text{denied the})$

2.5 allegations

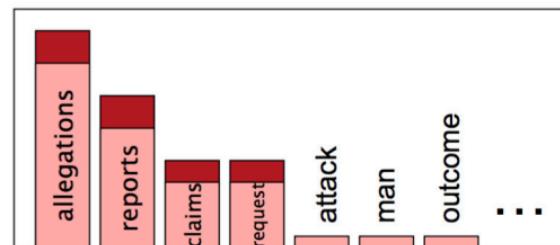
1.5 reports

0.5 claims

0.5 request

2 other

7 total

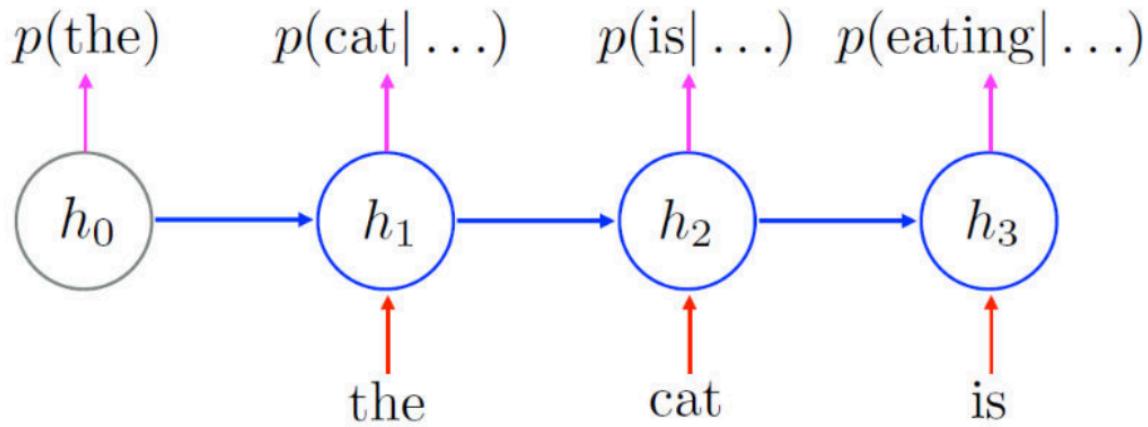


- Smoothing: 概率为 0 的问题

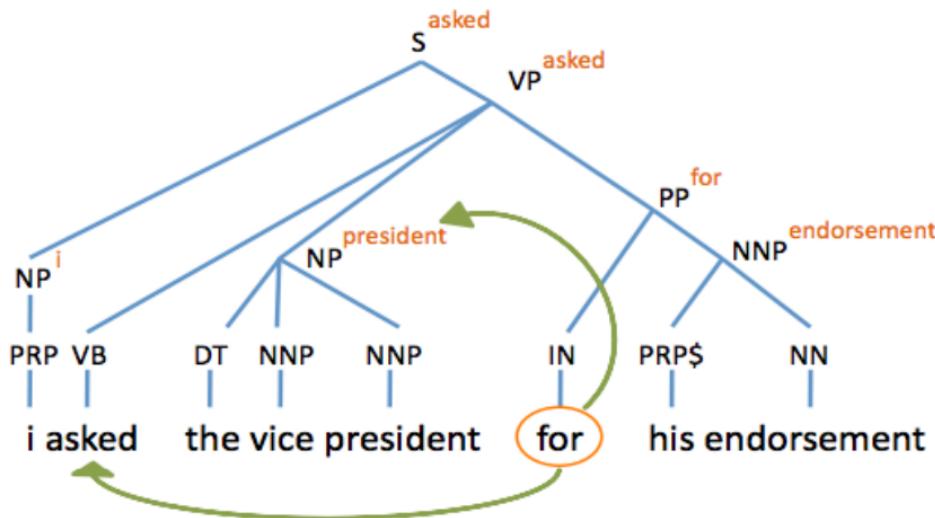
- ① Laplace (add-one, add- $\alpha$ ) [unseen words same probability]
- ② Deleted (held-out) estimation [things-we-saw-once]
- ③ Good Turing [things new, assuming leave-one-out]
- ④ Back-off (Katz) smoothing [non-linear zero-count]
- ⑤ Interpolation (Jelinek-Mercer) smoothing [linear zero-count]
- ⑥ Witten-Bell [interpolation, diversity of predicted word]
- ⑦ Absolute discounting [approximating Good-Turing]
- ⑧ Kneser-Ney [back-off, diversity of history]
- ⑨ Modified Kneser-Ney [interpolated back-off]

- 马尔可夫假设：N-Gram 模型

$p(\text{the}, \text{cat}, \text{is}, \text{eating})$



- 长距离依赖问题 Long-Term Dependency



- Log-Linear 语言模型

$$\hat{p}(y | x) \stackrel{\text{def}}{=} \frac{\text{count}(x, y)}{\text{count}(x)}$$



$$\begin{aligned} u(x, y) &\stackrel{\text{def}}{=} \exp \sum_{k=1}^K (\theta_k \cdot f_k(x, y)) \\ &= \exp(\vec{\theta} \cdot \vec{f}(x, y)) > 0 \end{aligned}$$

$$\hat{p}(y | x) \stackrel{\text{def}}{=} \frac{u(x, y)}{Z(x)}$$

$$p(y | x) \stackrel{\text{def}}{=} \frac{1}{Z(x)} \exp \sum_{k=1}^K \theta_k \cdot f_k(x, y)$$

$$\begin{aligned} Z(x) &\stackrel{\text{def}}{=} \sum_y u(x, y) \\ \sum_y p(y | x) &= 1 \end{aligned}$$

$$p(y | x) \stackrel{\text{def}}{=} \frac{1}{Z(x)} \prod_{\substack{k \text{ such that} \\ f_k(x, y) = 1}} (\exp \theta_k)$$

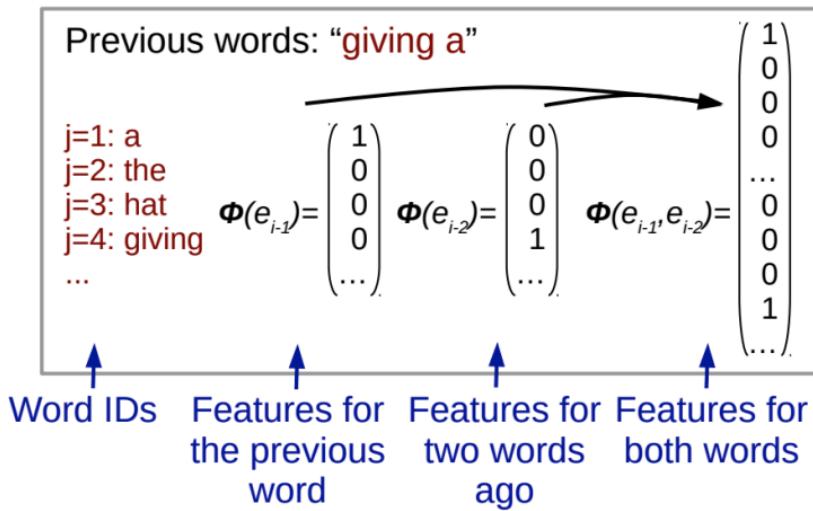
$$p(\text{solid circle}) \stackrel{\text{def}}{=} \frac{\exp(\theta_{\text{circle}} + \theta_{\text{solid}})}{\exp(\theta_{\text{circle}} + \theta_{\text{solid}}) + \exp(\theta_{\text{circle}} + \theta_{\text{striped}}) + \exp(\theta_{\text{triangle}} + \theta_{\text{solid}}) + \exp(\theta_{\text{triangle}} + \theta_{\text{striped}})}$$

$$p(\text{striped circle}) \stackrel{\text{def}}{=} \frac{\exp(\theta_{\text{circle}} + \theta_{\text{striped}})}{\exp(\theta_{\text{circle}} + \theta_{\text{solid}}) + \exp(\theta_{\text{circle}} + \theta_{\text{striped}}) + \exp(\theta_{\text{triangle}} + \theta_{\text{solid}}) + \exp(\theta_{\text{triangle}} + \theta_{\text{striped}})}$$

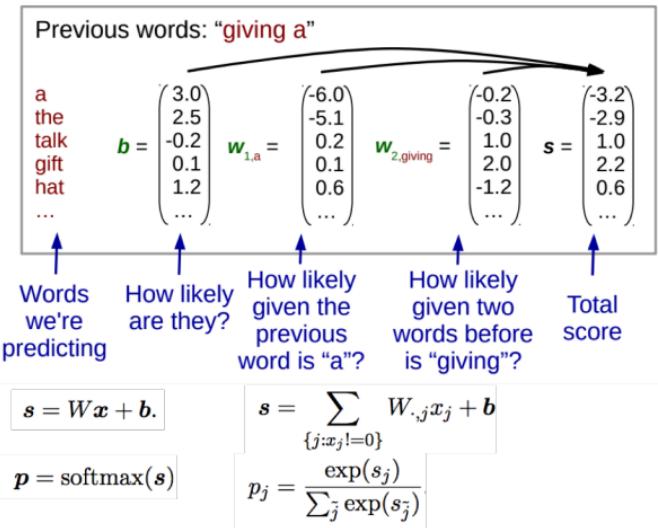
$$p(\text{solid triangle}) \stackrel{\text{def}}{=} \frac{\exp(\theta_{\text{triangle}} + \theta_{\text{solid}})}{\exp(\theta_{\text{circle}} + \theta_{\text{solid}}) + \exp(\theta_{\text{circle}} + \theta_{\text{striped}}) + \exp(\theta_{\text{triangle}} + \theta_{\text{solid}}) + \exp(\theta_{\text{triangle}} + \theta_{\text{striped}})}$$

$$p(\text{striped triangle}) \stackrel{\text{def}}{=} \frac{\exp(\theta_{\text{triangle}} + \theta_{\text{striped}})}{\exp(\theta_{\text{circle}} + \theta_{\text{solid}}) + \exp(\theta_{\text{circle}} + \theta_{\text{striped}}) + \exp(\theta_{\text{triangle}} + \theta_{\text{solid}}) + \exp(\theta_{\text{triangle}} + \theta_{\text{striped}})}$$

- Log-Linear 特征向量



- Log-Linear 到 Softmax 激活



- Log-Linear 优化更新

$$\mathbf{x} = \phi(e_{t-m+1}^{t-1})$$

$$\mathbf{s} = \sum_{\{j: x_j \neq 0\}} W_{\cdot, j} x_j + \mathbf{b}$$

$$\mathbf{p} = \text{softmax}(\mathbf{s})$$

$$\ell = -\log \mathbf{p}_{e_t}.$$



$$\frac{d\ell(e_{t-n+1}^t, W, \mathbf{b})}{d\mathbf{b}} = \frac{d\ell}{d\mathbf{p}} \frac{d\mathbf{p}}{d\mathbf{s}} \frac{d\mathbf{s}}{d\mathbf{b}}$$

$$\frac{d\ell(e_{t-n+1}^t, W, \mathbf{b})}{dW_{\cdot, j}} = \frac{d\ell}{d\mathbf{p}} \frac{d\mathbf{p}}{d\mathbf{s}} \frac{d\mathbf{s}}{dW_{\cdot, j}}$$



$$\frac{d\ell(e_{t-n+1}^t, W, \mathbf{b})}{d\mathbf{b}} = \mathbf{p} - \text{onehot}(e_t)$$

$$\frac{d\ell(e_{t-n+1}^t, W, \mathbf{b})}{dW_{\cdot, j}} = x_j (\mathbf{p} - \text{onehot}(e_t))$$

- Log-Linear 优化更新推导

$$P(y \mid x, \mathbf{W}) = \frac{e^{\sum_k \mathbf{W}_k \phi_k(x, y)}}{\sum_{y' \in \mathcal{Y}} e^{\sum_k \mathbf{W}_k \phi_k(x, y')}}$$

$$\log P(y \mid x, \mathbf{W}) = \underbrace{\mathbf{W} \cdot \phi(x, y)}_{\text{Linear term}} - \underbrace{\log \sum_{y' \in \mathcal{Y}} e^{\mathbf{W} \cdot \phi(x, y')}}_{\text{Normalization term}}$$

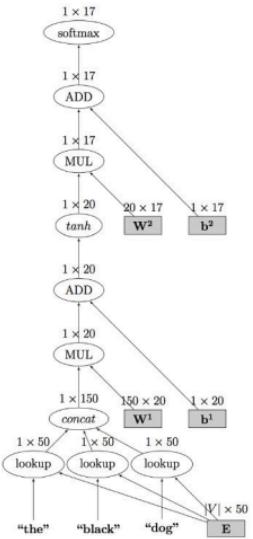
$$\begin{aligned} L(\mathbf{W}) &= \sum_{i=1}^n \log P(y_i \mid x_i) \\ &= \sum_{i=1}^n \mathbf{W} \cdot \phi(x_i, y_i) - \sum_{i=1}^n \log \sum_{y' \in \mathcal{Y}} e^{\mathbf{W} \cdot \phi(x_i, y')} \end{aligned}$$

$$\begin{aligned} \left. \frac{dL}{d\mathbf{W}} \right|_{\mathbf{W}} &= \sum_{i=1}^n \phi(x_i, y_i) - \sum_{i=1}^n \frac{\sum_{y' \in \mathcal{Y}} \phi(x_i, y') e^{\mathbf{W} \cdot \phi(x_i, y')}}{\sum_{z' \in \mathcal{Y}} e^{\mathbf{W} \cdot \phi(x_i, z')}} \\ &= \sum_{i=1}^n \phi(x_i, y_i) - \sum_{i=1}^n \sum_{y' \in \mathcal{Y}} \phi(x_i, y') \frac{e^{\mathbf{W} \cdot \phi(x_i, y')}}{\sum_{z' \in \mathcal{Y}} e^{\mathbf{W} \cdot \phi(x_i, z')}} \\ &= \underbrace{\sum_{i=1}^n \phi(x_i, y_i)}_{\text{Empirical counts}} - \underbrace{\sum_{i=1}^n \sum_{y' \in \mathcal{Y}} \phi(x_i, y') P(y' \mid x_i, \mathbf{W})}_{\text{Expected counts}} \end{aligned}$$

$$\Delta = \left. \frac{dL}{d\mathbf{W}} \right|_{\mathbf{W}}$$

$$\mathbf{W} \leftarrow \mathbf{W} + \beta_* \Delta$$

- NNLM: 加  $\tanh(x)$  截断的 Log-Linear

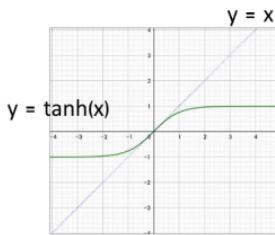
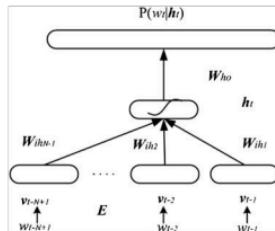


$$\mathbf{h}' = W_{xh} \mathbf{x} + b_h$$

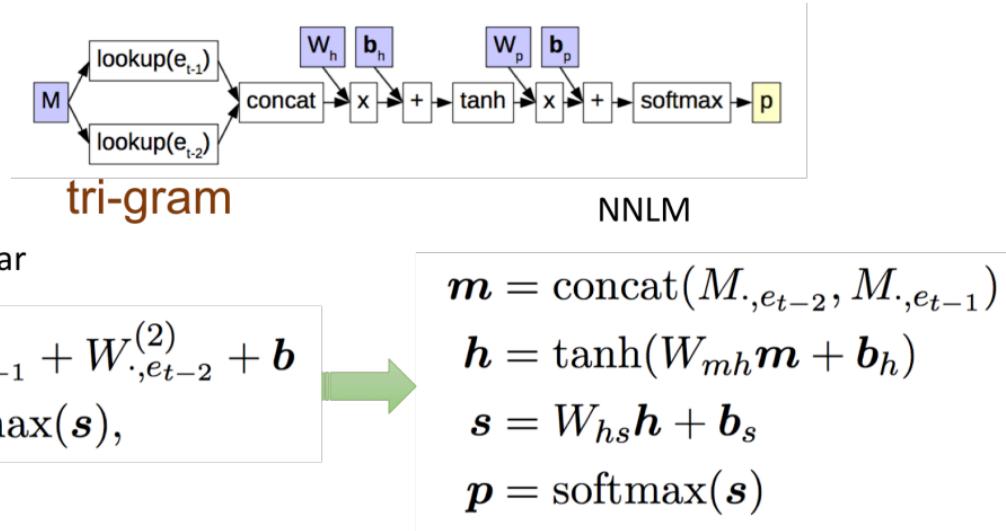
$$\mathbf{h} = \tanh(\mathbf{h}')$$

$$y = w_{hy} \mathbf{h} + b_y$$

$$\mathbf{p} = \text{softmax}(y)$$



- NNLM (trigram) : 从加和变成拼接



- NNLM 小结

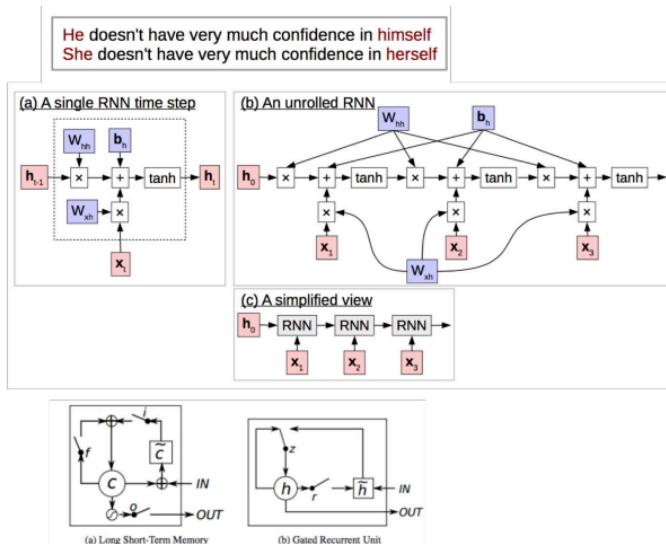
1. 吸收了 Log-Linear 的优点，可以跳跃周边词汇
  - 让上下文表达变得灵活
  - 更强的泛化能力
2. 加和变成拼接，减少信息丢失
  - 灵活的组合方式
3. 依然对 N 有限制
  - 无法彻底解决长距离依赖问题

## 4.2 从神经网络模型到 RNN 模型

- RNN-LM 模型

RNN  
顺序  
信息压缩

记忆学习



- RNN-LM 模型

## 上下文宽带限制

$$\begin{aligned}\mathbf{m} &= \text{concat}(M_{\cdot, e_{t-2}}, M_{\cdot, e_{t-1}}) \\ \mathbf{h} &= \tanh(W_{mh}\mathbf{m} + \mathbf{b}_h) \\ \mathbf{s} &= W_{hs}\mathbf{h} + \mathbf{b}_s \\ \mathbf{p} &= \text{softmax}(\mathbf{s})\end{aligned}$$



## 顺序压缩

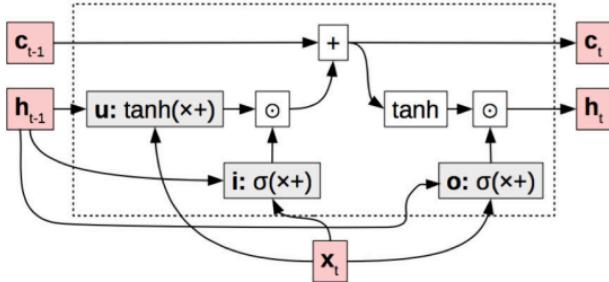
$$\begin{aligned}\mathbf{m}_t &= M_{\cdot, e_{t-1}} \\ \mathbf{h}_t &= \begin{cases} \tanh(W_{mh}\mathbf{m}_t + W_{hh}\mathbf{h}_{t-1} + \mathbf{b}_h) & t \geq 1, \\ \mathbf{0} & \text{otherwise.} \end{cases} \\ \mathbf{p}_t &= \text{softmax}(W_{hs}\mathbf{h}_t + \mathbf{b}_s).\end{aligned}$$



## 带记忆顺序压缩

$$\begin{aligned}\mathbf{m}_t &= M_{\cdot, e_{t-1}} \\ \mathbf{h}_t &= \text{RNN}(\mathbf{m}_t, \mathbf{h}_{t-1}) \\ \mathbf{p}_t &= \text{softmax}(W_{hs}\mathbf{h}_t + \mathbf{b}_s).\end{aligned}$$

- RNN-LM 模型: LSTM 和 GRU



update  $u$ : what value do we try to add to the memory cell?  
 input  $i$ : how much of the update do we allow to go through?  
 output  $o$ : how much of the cell do we reflect in the next state?

$$\sigma(x) = \frac{1}{1 + \exp(-x)}$$

$$z = x \odot y$$

$$z_i = x_i * y_i$$

## LSTM

$$\begin{aligned} u_t &= \tanh(W_{xu}x_t + W_{hu}h_{t-1} + b_u) \\ i_t &= \sigma(W_{xi}x_t + W_{hi}h_{t-1} + b_i) \\ f_t &= \sigma(W_{xf}x_t + W_{hf}h_{t-1} + b_f) \\ o_t &= \sigma(W_{xo}x_t + W_{ho}h_{t-1} + b_o) \\ c_t &= i_t \odot u_t + f_t \odot c_{t-1} \\ h_t &= o_t \odot \tanh(c_t). \end{aligned}$$

## GRU

$$\begin{aligned} r_t &= \sigma(W_{xr}x_t + W_{hr}h_{t-1} + b_r) \\ z_t &= \sigma(W_{xz}x_t + W_{hz}h_{t-1} + b_z) \\ \tilde{h}_t &= \tanh(W_{xh}x_t + W_{hh}(r_t \odot h_{t-1}) + b_h) \\ h_t &= (1 - z_t)h_{t-1} + z_t\tilde{h}_t. \end{aligned}$$

- RNN-LM 模型：批处理

---

**Algorithm 2** A batch learning algorithm

---

```
1: procedure BATCH
2:   for several epochs of training do
3:     for each training example in the data do
4:       Calculate and accumulate gradients of the loss
5:     end for
6:     Update the parameters according to the accumulated gradient
7:   end for
8: end procedure
```

---

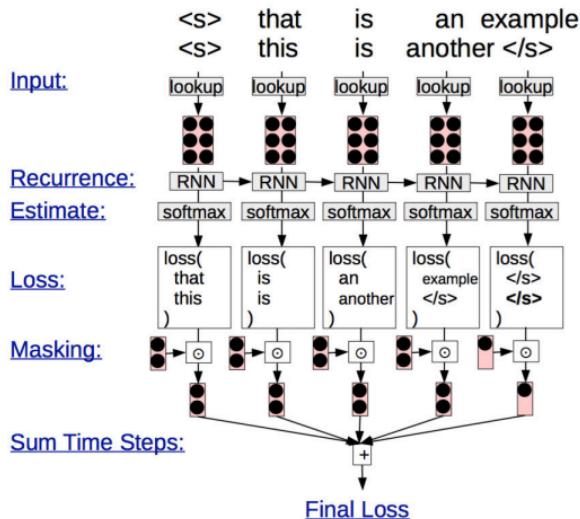
Operations w/o Minibatching

$$\tanh(\begin{matrix} W & x_1 & b \end{matrix} + \begin{matrix} W & x_2 & b \end{matrix} + \begin{matrix} W & x_3 & b \end{matrix})$$

Operations with Minibatching

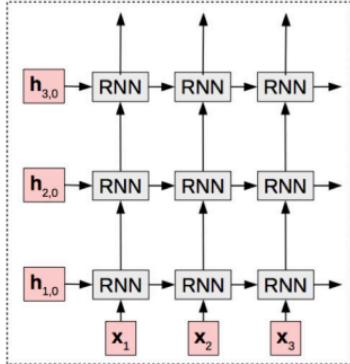
$$\begin{matrix} x_1 & x_2 & x_3 \end{matrix} \xrightarrow{\text{concat}} \begin{matrix} W & X & B \end{matrix} \xrightarrow{\text{broadcast}} b$$
$$\tanh(\begin{matrix} W & X & B \end{matrix} + \begin{matrix} W & X & B \end{matrix})$$

- RNN-LM 模型：批处理

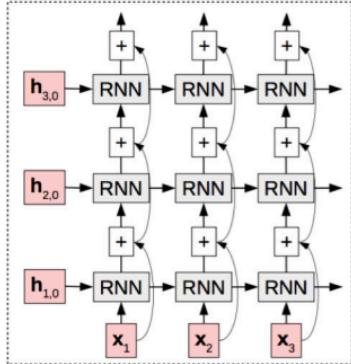


- 扩展 RNN-LM 模型：堆栈和残差

(a) A stacked RNN



(b) With residual connections



$$\mathbf{h}_{1,t} = \text{RNN}_1(\mathbf{x}_t, \mathbf{h}_{1,t-1})$$

$$\mathbf{h}_{2,t} = \text{RNN}_2(\mathbf{h}_{1,t}, \mathbf{h}_{2,t-1})$$

$$\mathbf{h}_{3,t} = \text{RNN}_3(\mathbf{h}_{2,t}, \mathbf{h}_{3,t-1}),$$

$$\mathbf{h}_{1,t} = \text{RNN}_1(\mathbf{x}_t, \mathbf{h}_{1,t-1}) + \mathbf{x}_t$$

$$\mathbf{h}_{2,t} = \text{RNN}_2(\mathbf{h}_{1,t}, \mathbf{h}_{2,t-1}) + \mathbf{h}_{1,t}$$

$$\mathbf{h}_{3,t} = \text{RNN}_3(\mathbf{h}_{2,t}, \mathbf{h}_{3,t-1}) + \mathbf{h}_{2,t}.$$

- RNN-LM 小结

1. RNN 解决长距离依赖问题

- 保留了 NNLM 的优点：上下文、拼接等

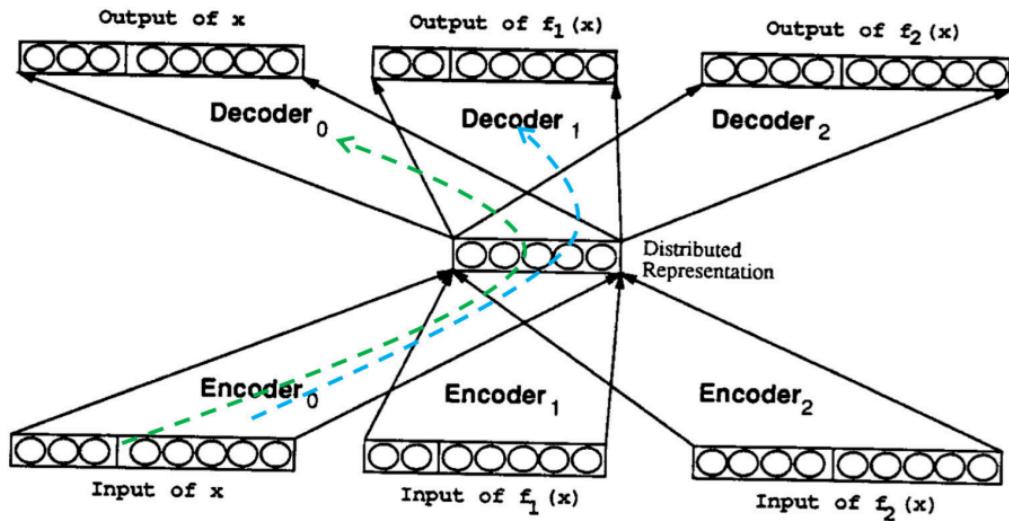
2. RNN 变种的各种优点

- 记忆能力
- 多种顺序信息压缩

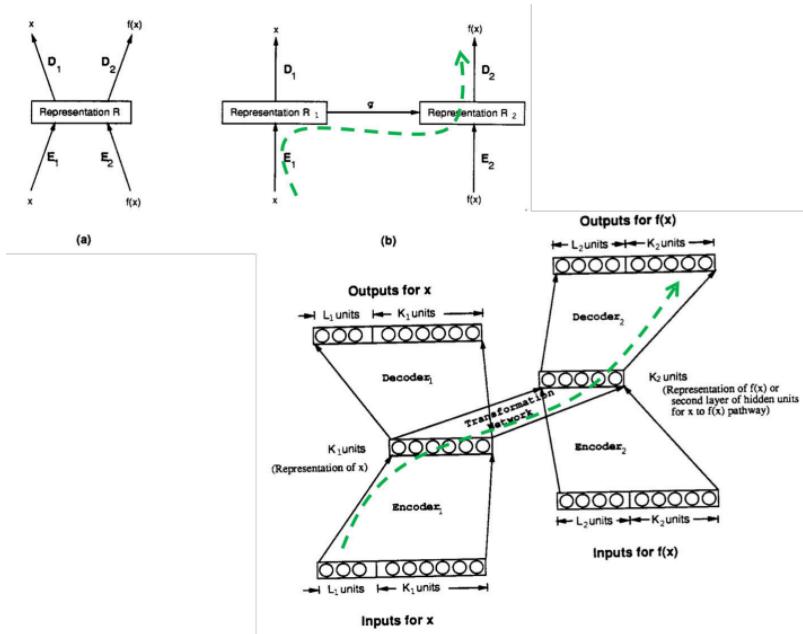
3. 问题：增大了计算量

### 4.3 从 RNN 模型到 Encoder-Decoder 模型

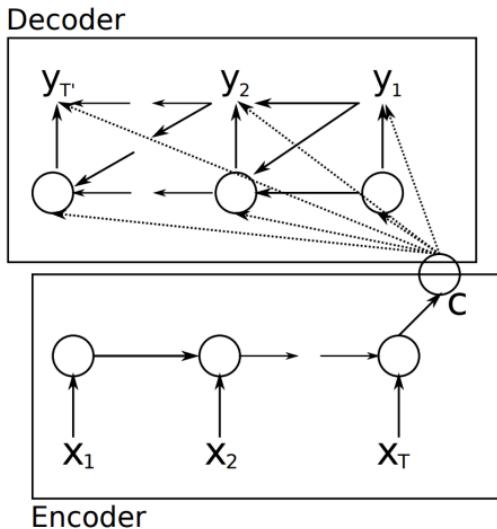
- 端到端翻译：双口 Encoder-Decoder 模型



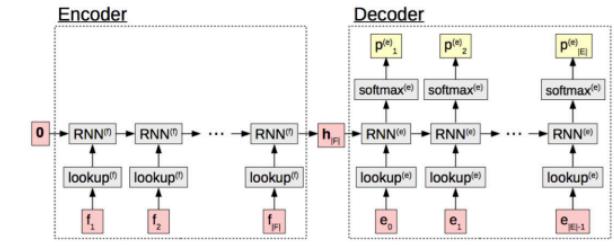
- 端到端翻译：带转换的双口 Encoder-Decoder 模型



- 端到端翻译：Encoder-Decoder 模型 + RNN 语言模型



- RNN Encoder-Decoder 模型



$$\mathbf{m}_t^{(f)} = M_{,f_t}^{(f)}$$

$$\mathbf{h}_t^{(f)} = \begin{cases} \text{RNN}^{(f)}(\mathbf{m}_t^{(f)}, \mathbf{h}_{t-1}^{(f)}) & t \geq 1, \\ \mathbf{0} & \text{otherwise.} \end{cases}$$

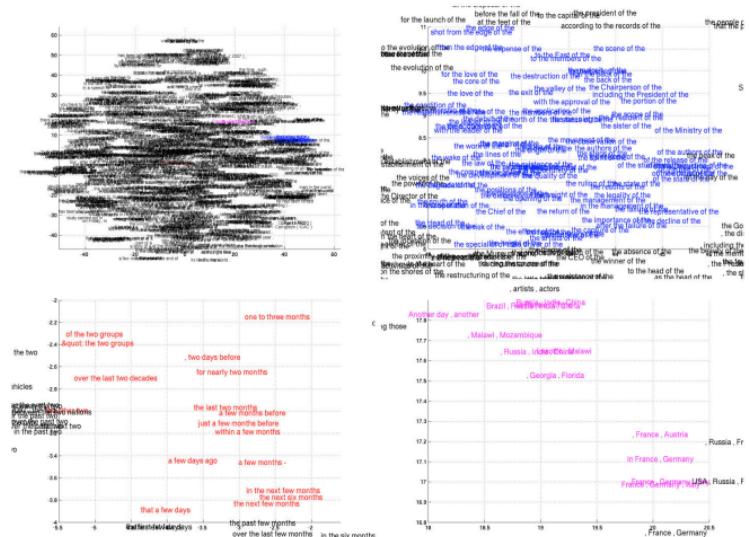
$$\mathbf{m}_t^{(e)} = M_{,e_{t-1}}^{(e)}$$

$$\mathbf{h}_t^{(e)} = \begin{cases} \text{RNN}^{(e)}(\mathbf{m}_t^{(e)}, \mathbf{h}_{t-1}^{(e)}) & t \geq 1, \\ \mathbf{h}_{|F|}^{(f)} & \text{otherwise.} \end{cases}$$

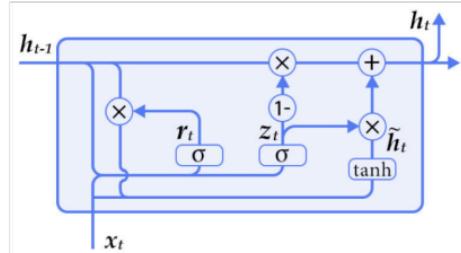
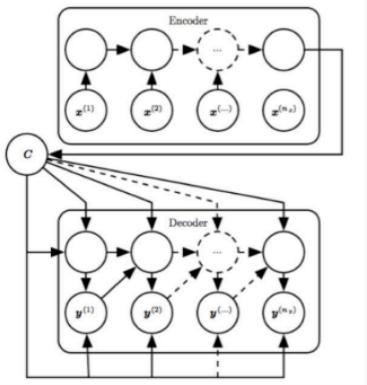
$$\mathbf{p}_t^{(e)} = \text{softmax}(W_{hs}\mathbf{h}_t^{(e)} + b_s)$$



- RNN Encoder-Decoder 自带词嵌入效果



- GRU Encoder–Decoder



$$\mathbf{c} = \tanh(\mathbf{V}\mathbf{h}^{(N)})$$

$$h_j^{(t)} = z_j h_j^{(t-1)} + (1 - z_j) \tilde{h}_j^{(t)}$$

$$\tilde{h}_j^{(t)} = \tanh([\mathbf{W}_e(\mathbf{x}_t)]_j + [\mathbf{U}(\mathbf{r} \odot \mathbf{h}_{(t-1)})]_j),$$

$$z_j = \sigma([\mathbf{W}_z e(\mathbf{x}_t)]_j + [\mathbf{U}_z \mathbf{h}_{(t-1)}]_j),$$

$$r_j = \sigma([\mathbf{W}_r e(\mathbf{x}_t)]_j + [\mathbf{U}_r \mathbf{h}_{(t-1)}]_j).$$

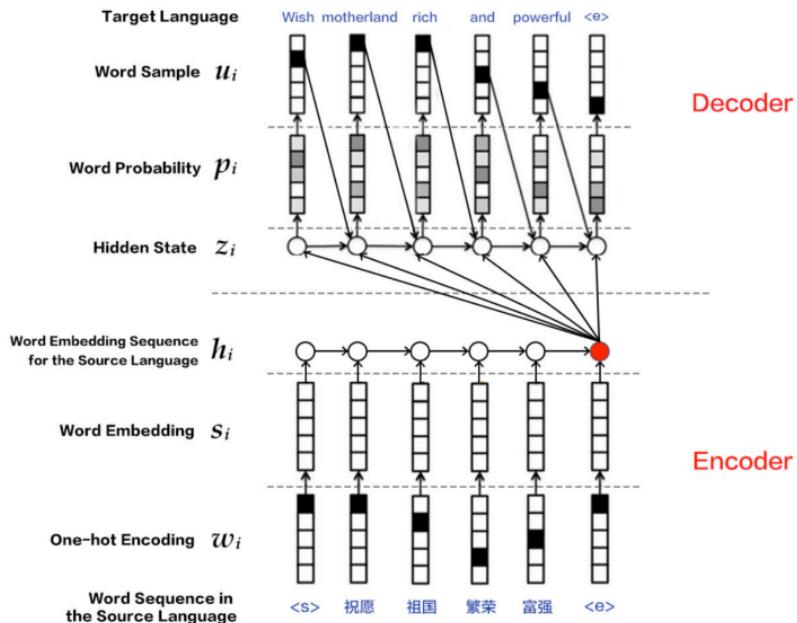
$$h'_j^{(t)} = z'_j h'_j^{(t-1)} + (1 - z'_j) \tilde{h}'_j^{(t)}$$

$$\tilde{h}'_j^{(t)} = \tanh([\mathbf{W}'_e(\mathbf{y}_{t-1})]_j + r'_j [\mathbf{U}' \mathbf{h}'_{(t-1)} + \mathbf{C} \mathbf{c}])$$

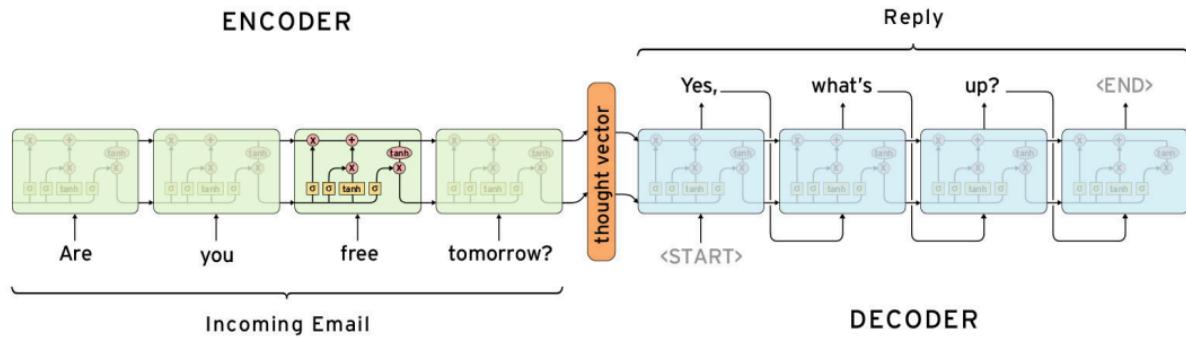
$$z'_j = \sigma([\mathbf{W}'_z e(\mathbf{y}_{t-1})]_j + [\mathbf{U}'_z \mathbf{h}'_{(t-1)}]_j + [\mathbf{C}_z \mathbf{c}]_j)$$

$$r'_j = \sigma([\mathbf{W}'_r e(\mathbf{y}_{t-1})]_j + [\mathbf{U}'_r \mathbf{h}'_{(t-1)}]_j + [\mathbf{C}_r \mathbf{c}]_j)$$

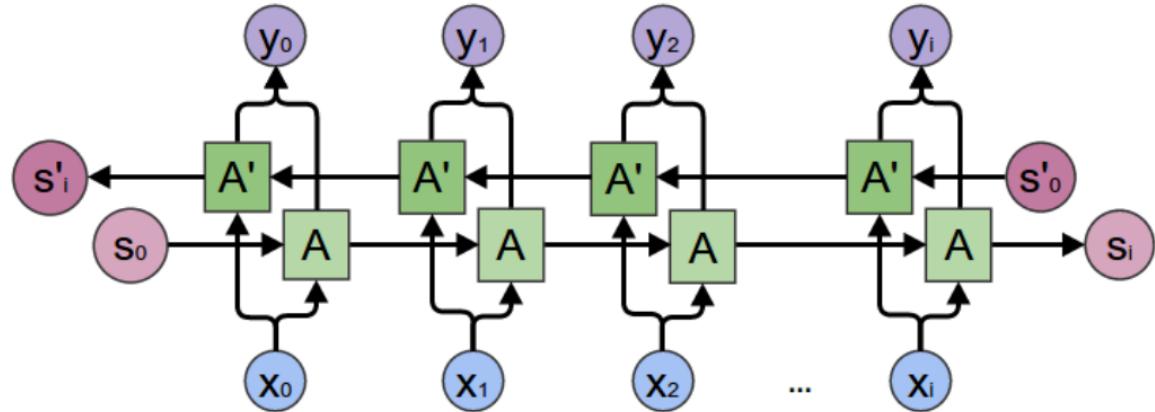
- RNN Encoder–Decoder: 端到端翻译



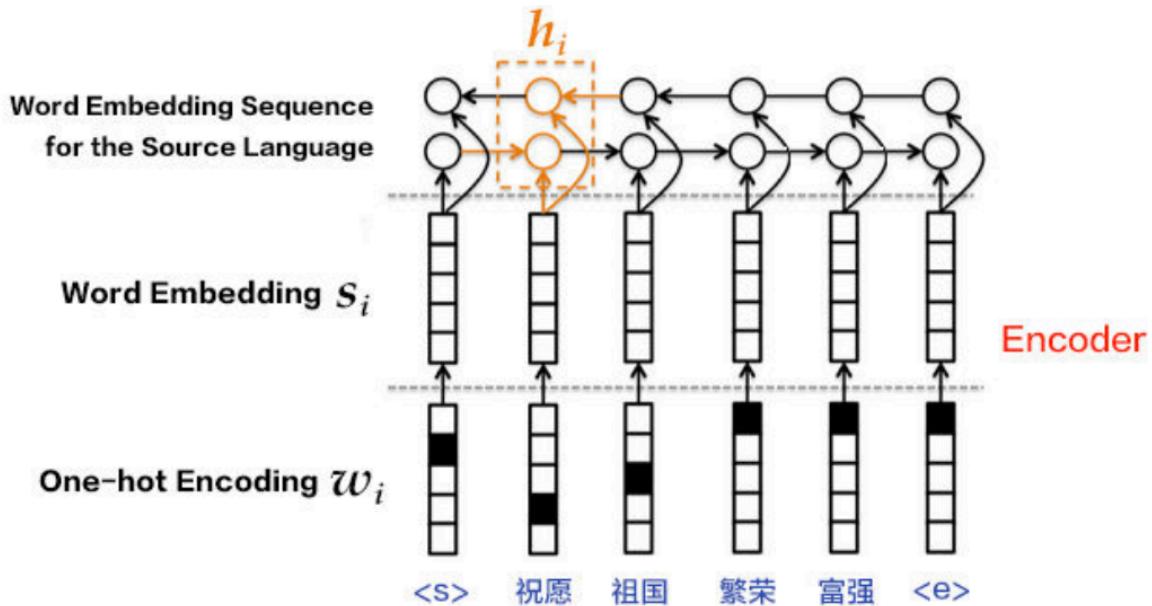
- LSTM Encoder–Decoder: 端到端问答！



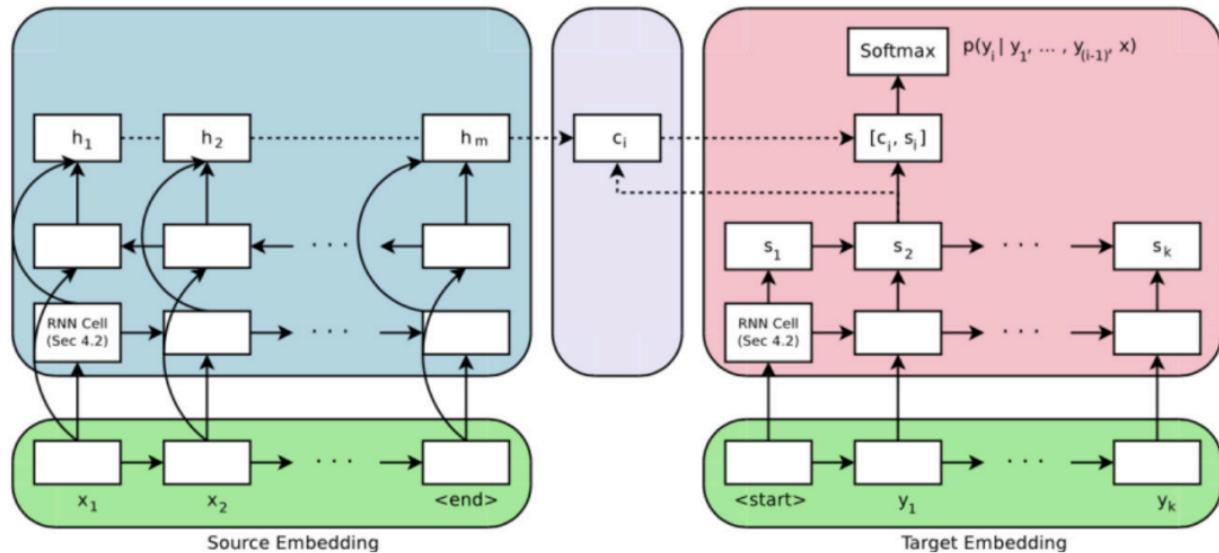
- 双向 RNN



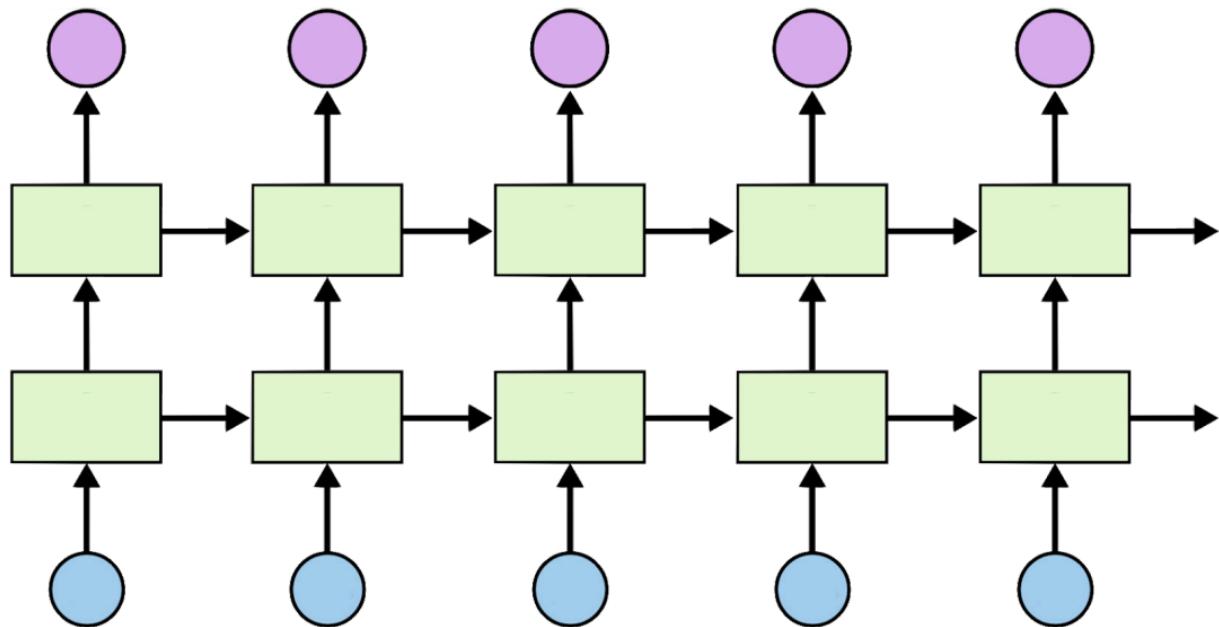
- 双向 RNN 词嵌入



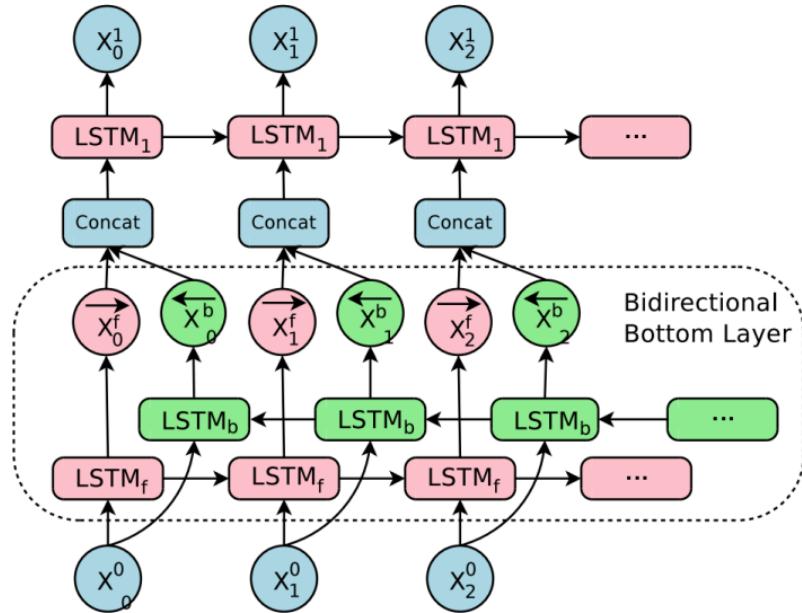
- 双向 RNN Encoder–Decoder



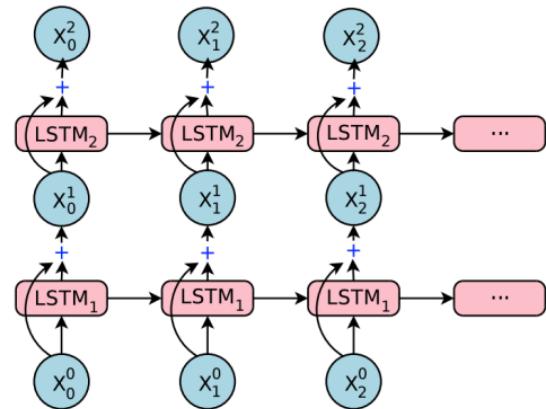
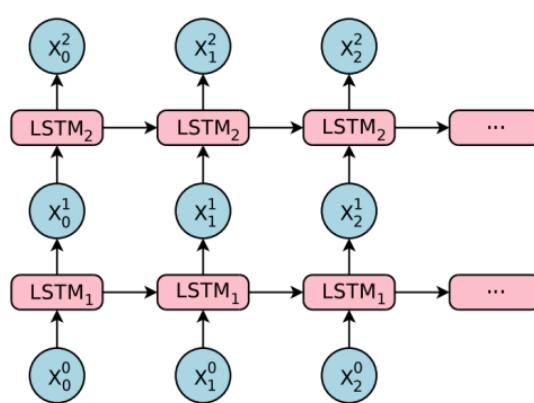
- 堆栈 RNN



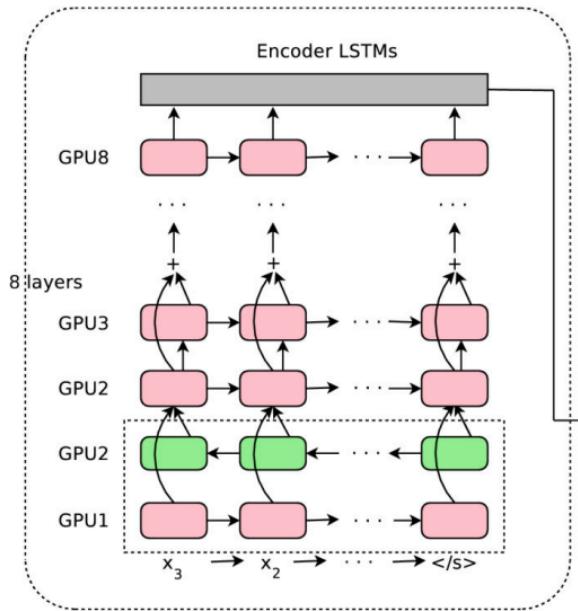
- 双向堆栈 RNN Encoder



- 带残差的堆栈 RNN



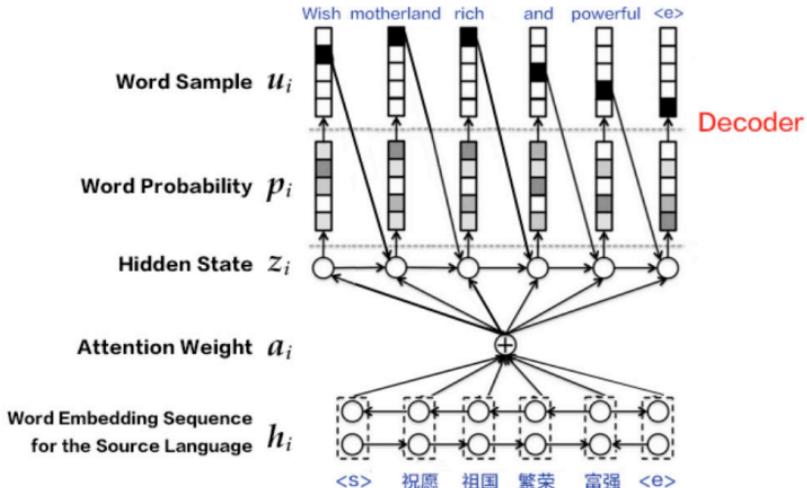
- GNMT 谷歌神经翻译的 Encoder



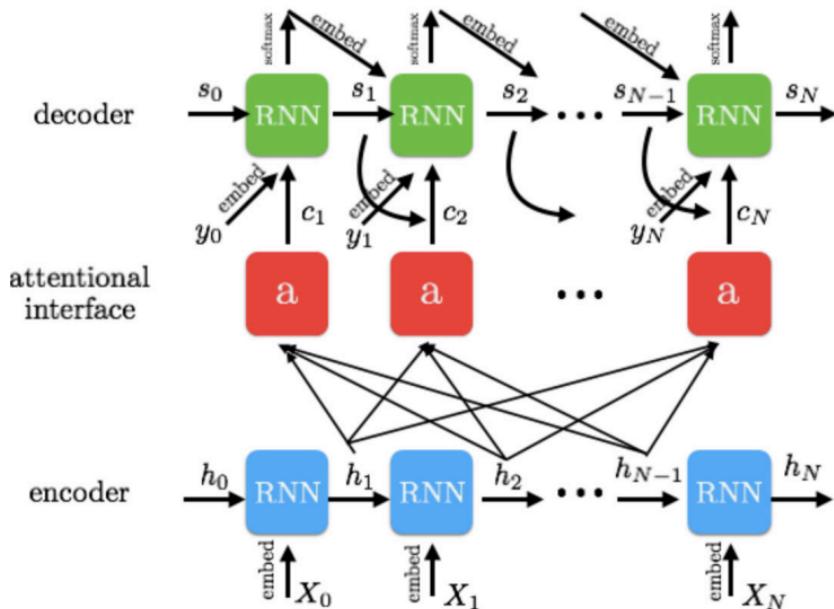
- RNN Encoder–Decoder 小结
  - 1. 使用 RNN 作为编解码网络
    - 双向 RNN
    - 双向 RNN + 堆栈 + 残差
    - 双向 RNN + 堆栈 + 残差
  - 2. 计算量增大，GPU 加速
  - 3. 原始的 Sequence to Sequence 模型

## 4.4 Attention 模型

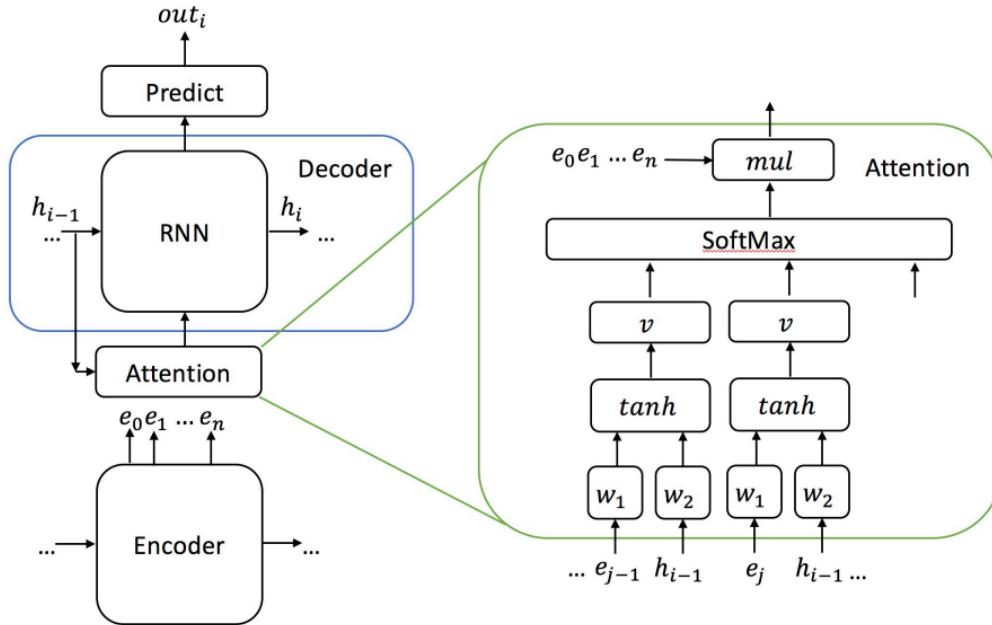
- 注意力机制的端到端翻译



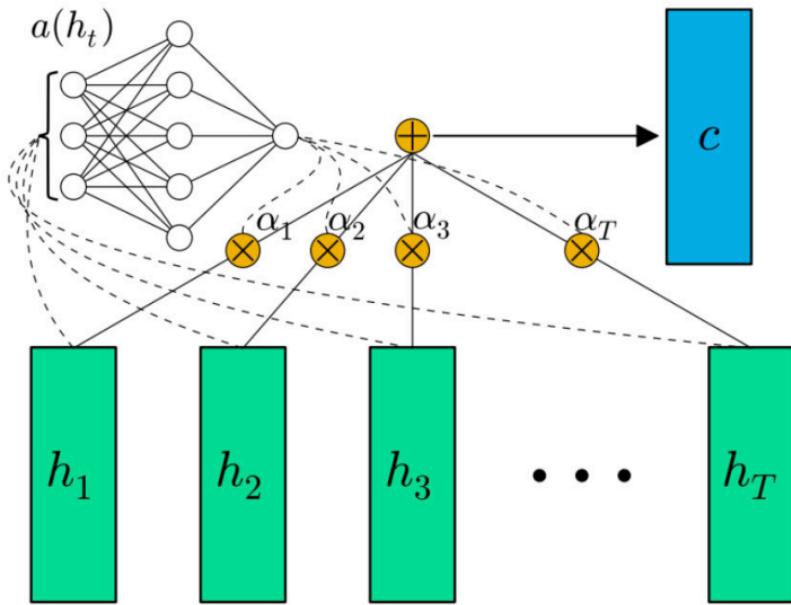
- 注意力机制的 Encoder-Decoder



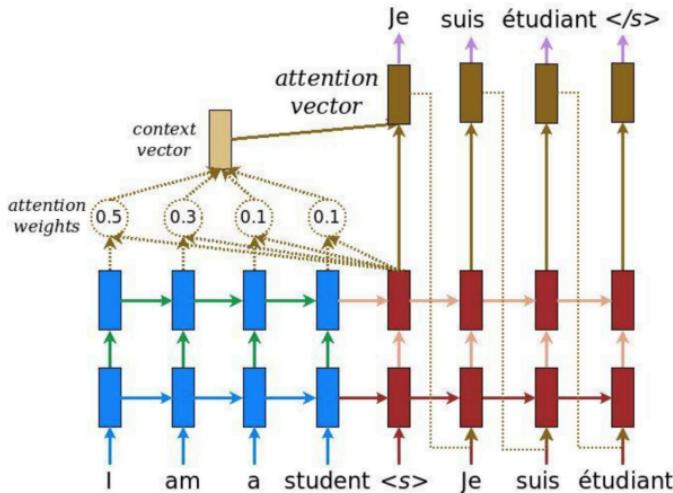
- 注意力网络



- 注意力权重



- 注意力机制：自动调整映射范围



- 注意力机制构成

### Attention

Weight

$$\alpha_{ts} = \frac{\exp(\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_s))}{\sum_{s'=1}^S \exp(\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_{s'}))}$$

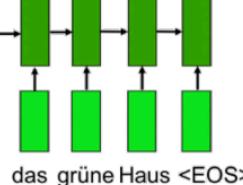
Context  
Vector

$$\mathbf{c}_t = \sum_s \alpha_{ts} \bar{\mathbf{h}}_s$$

$$\mathbf{a}_t = f(\mathbf{c}_t, \mathbf{h}_t) = \tanh(W_c[\mathbf{c}_t; \mathbf{h}_t])$$

Attention  
Vector

Encoder



Context Vector

Attention

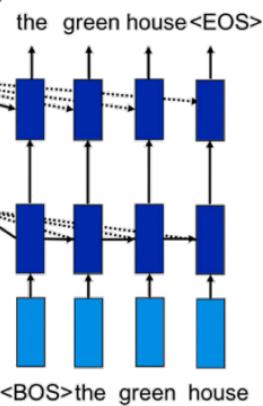
Attention Weight

Attention

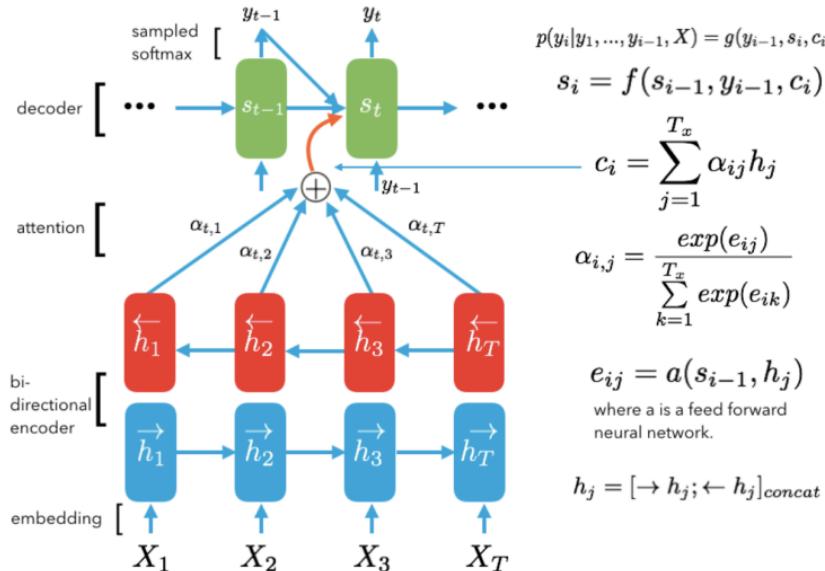
Vector

context

Decoder



- 注意力机制

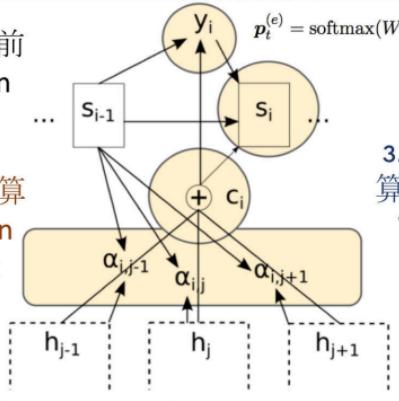


- 注意力机制计算流程

1. 根据Attention Score选择候选词

0. 已知当前  
Attention  
Score

2. 可以计算  
Attention  
Weight



$$\mathbf{h}_j^{(f)} = [\overleftarrow{\mathbf{h}}_j^{(f)}; \overrightarrow{\mathbf{h}}_j^{(f)}]$$

$$\overrightarrow{\mathbf{h}}_j^{(f)} = \text{RNN}(\text{embed}(f_j), \overrightarrow{\mathbf{h}}_{j-1}^{(f)})$$

$$\overleftarrow{\mathbf{h}}_j^{(f)} = \text{RNN}(\text{embed}(f_j), \overleftarrow{\mathbf{h}}_{j+1}^{(f)}).$$

3. 然后计  
算Context  
Vector

5. 根据候选词、  
Attention Vector  
计算Score

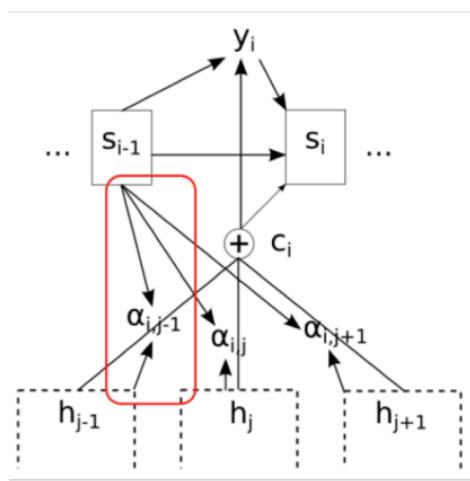
4. 再计算  
Attention  
Vector

$$\text{attn\_score}(\mathbf{h}_j^{(f)}, \mathbf{h}_t^{(e)}) := \mathbf{h}_j^{(f)\top} \mathbf{h}_t^{(e)}$$

$$\text{attn\_score}(\mathbf{h}_j^{(f)}, \mathbf{h}_t^{(e)}) := \mathbf{h}_j^{(f)\top} W_a \mathbf{h}_t^{(e)},$$

$$\text{attn\_score}(\mathbf{h}_t^{(e)}, \mathbf{h}_j^{(f)}) := \mathbf{w}_{a2}^\top \tanh(W_{a1}[\mathbf{h}_t^{(e)}, \mathbf{h}_j^{(f)}])$$

- 注意力得分 3 种常见计算方法



### 1. 简单的点积计算 (类似cosine距离)

$$\text{attn\_score}(\mathbf{h}_j^{(f)}, \mathbf{h}_t^{(e)}) := \mathbf{h}_j^{(f)\top} \mathbf{h}_t^{(e)}$$

### 2. 通过矩阵映射计算 (加权的cosine距离)

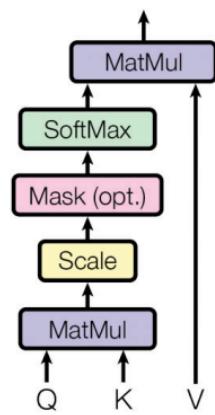
$$\text{attn\_score}(\mathbf{h}_j^{(f)}, \mathbf{h}_t^{(e)}) := \mathbf{h}_j^{(f)\top} W_a \mathbf{h}_t^{(e)}.$$

### 3. 通过神经网络计算

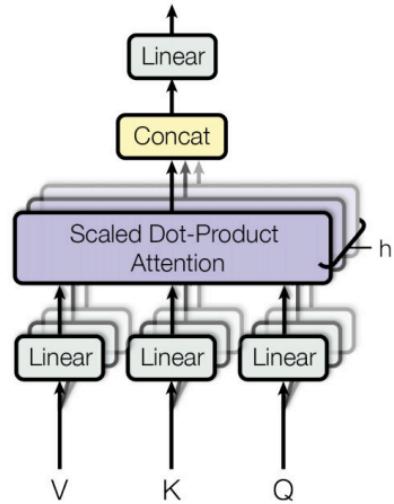
$$\text{attn\_score}(\mathbf{h}_t^{(e)}, \mathbf{h}_j^{(f)}) := \mathbf{w}_{a2}^\top \tanh(W_{a1}[\mathbf{h}_t^{(e)}; \mathbf{h}_j^{(f)}])$$

- 注意力得分网络

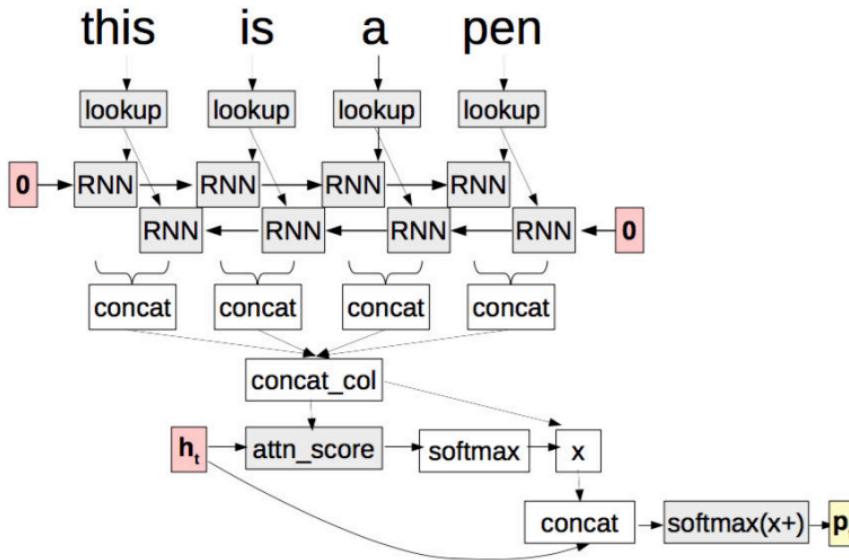
Scaled Dot-Product Attention



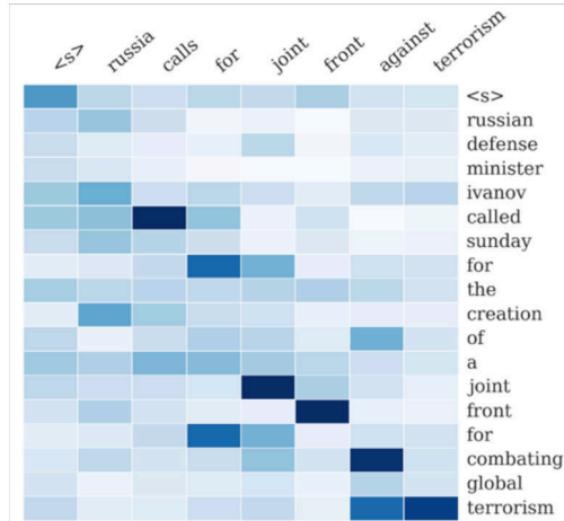
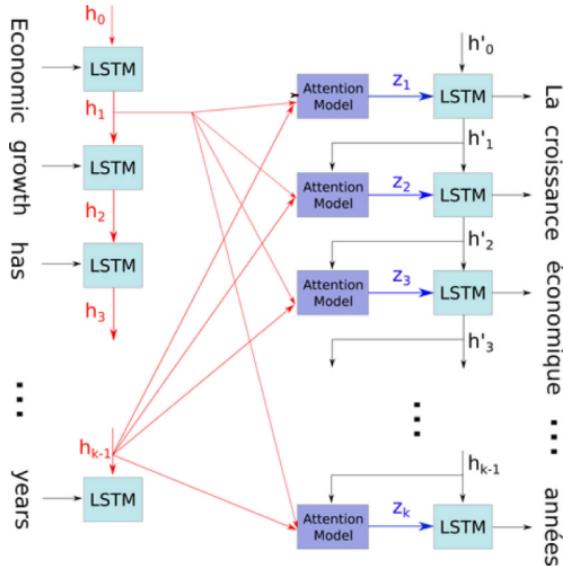
Multi-Head Attention



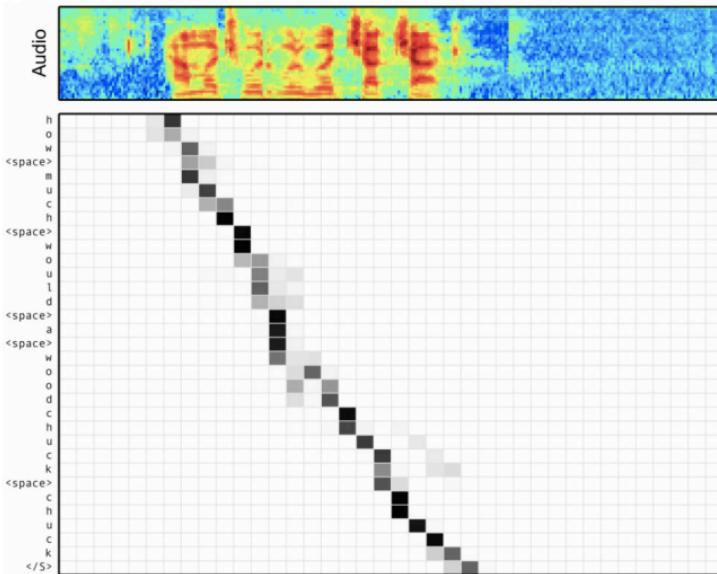
- 注意力计算图



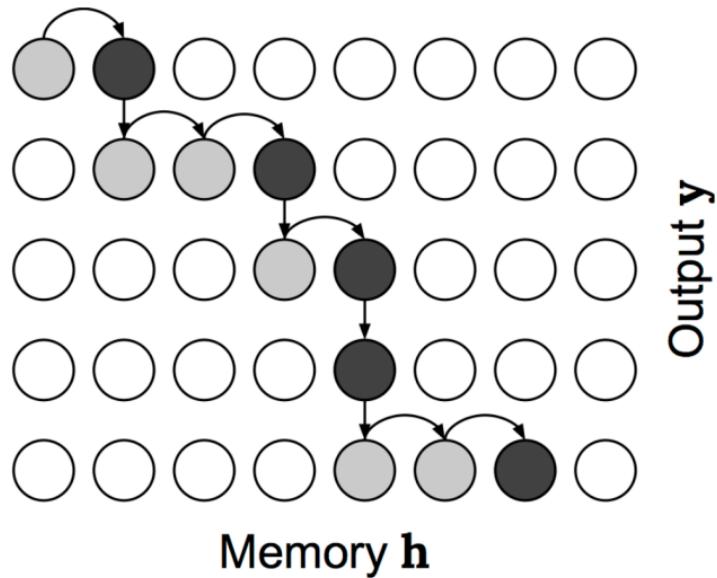
- 注意力得分和对齐



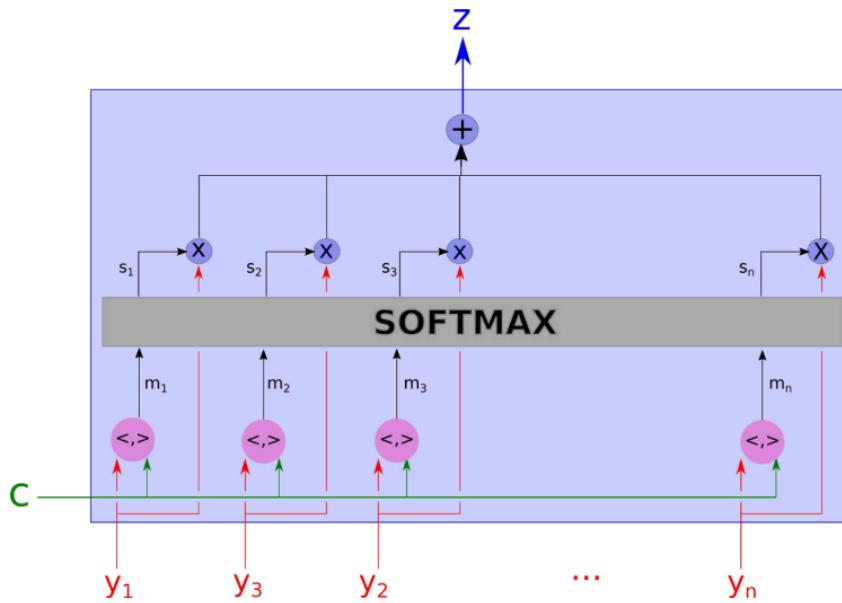
- 注意力机制对齐：语音和文字信号



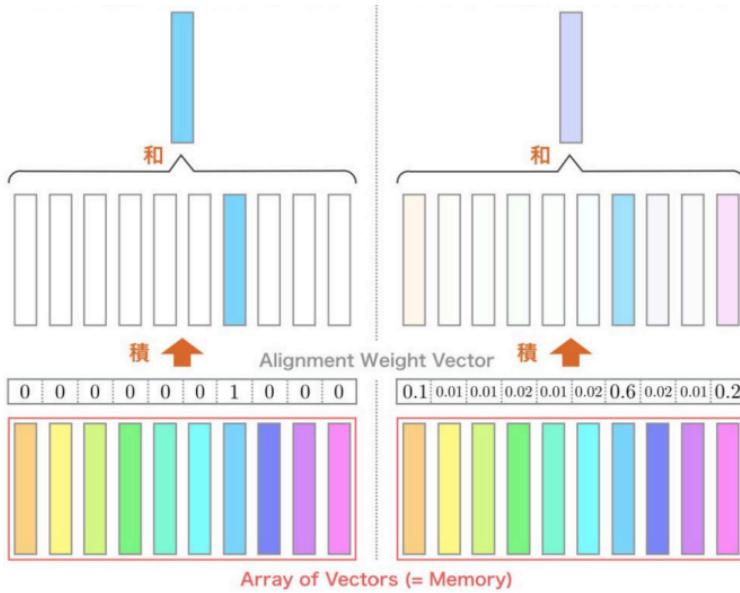
- 对齐机制



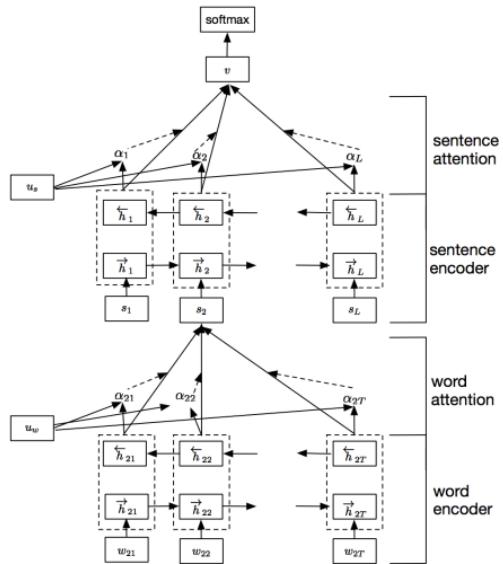
- 软注意力机制



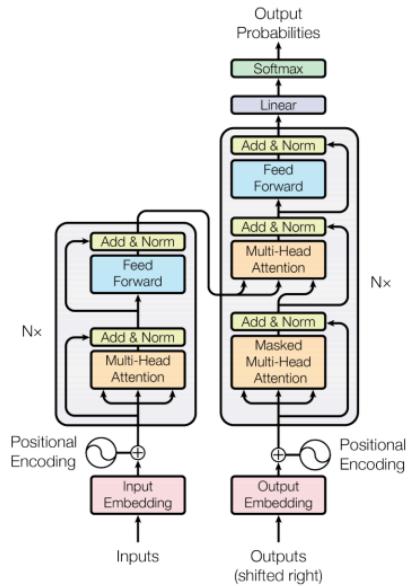
- 硬和软的注意力机制对比: 0/1 vs probability



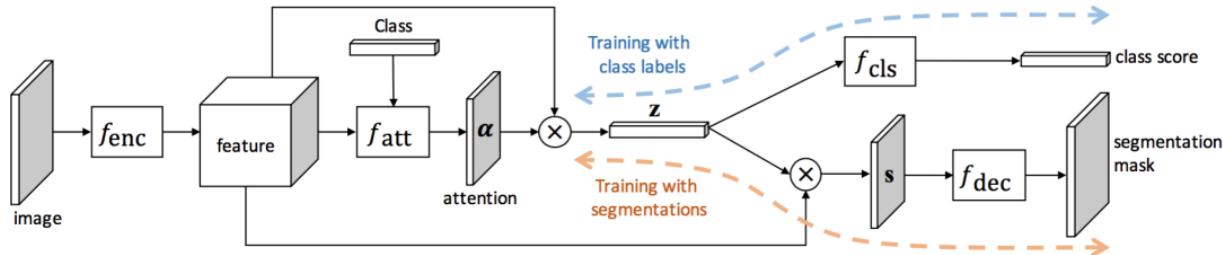
- 分层注意力机制



- Attention Is All You Need: Transformer Model



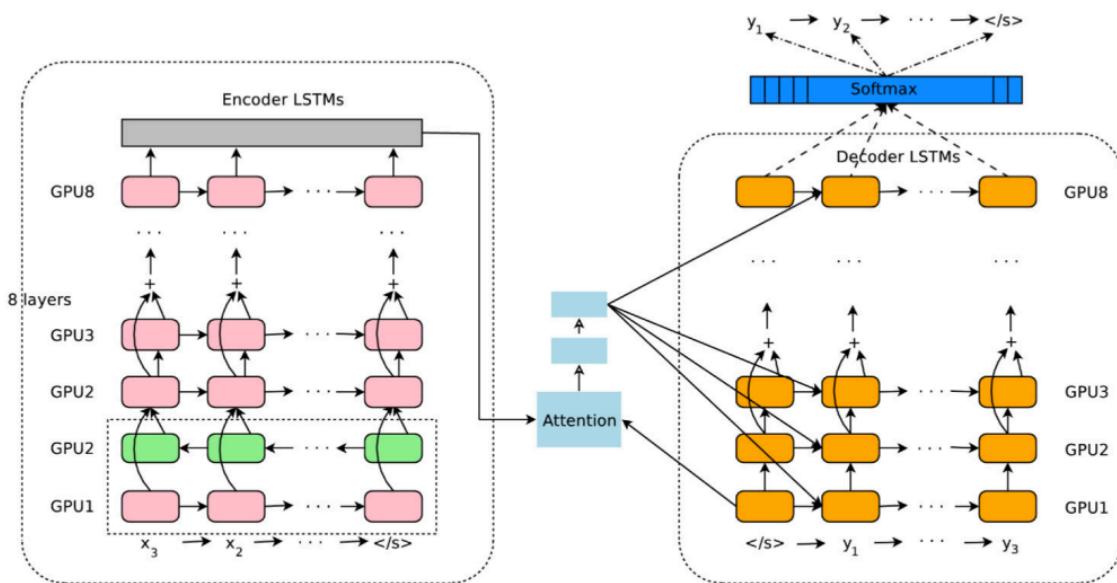
- 图像注意力机制



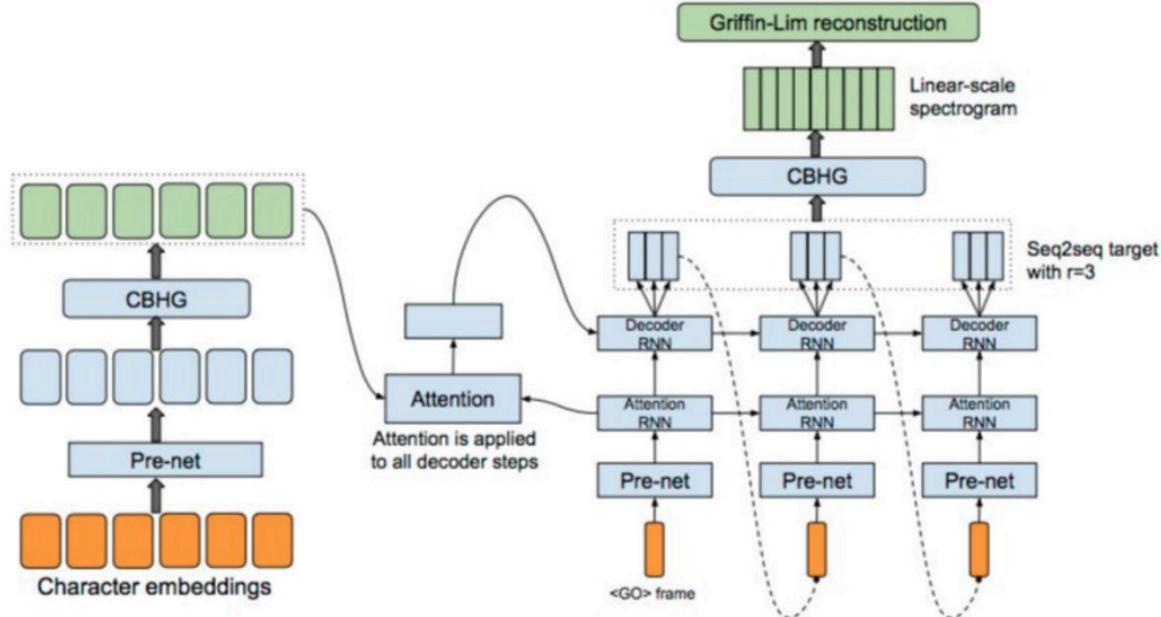
- Attention 机制小结
  1. Attention 是上下文环境的自适应强化
  2. Attention 可以有软硬的区别
  3. Attention 和传统的对齐有很强的关联
  4. Attention 计算分为 5 个步骤
    - 根据得分选定候选词
    - 计算权重
    - 计算上下文向量
    - 计算注意力向量
    - 重新计算得分
  5. Attention 得分主要有 3 种方法：点积、矩阵积、神经网络
  6. Attention 扩展：分层模型、转换器模型

## 4.5 Seq2Seq 模型：端到端

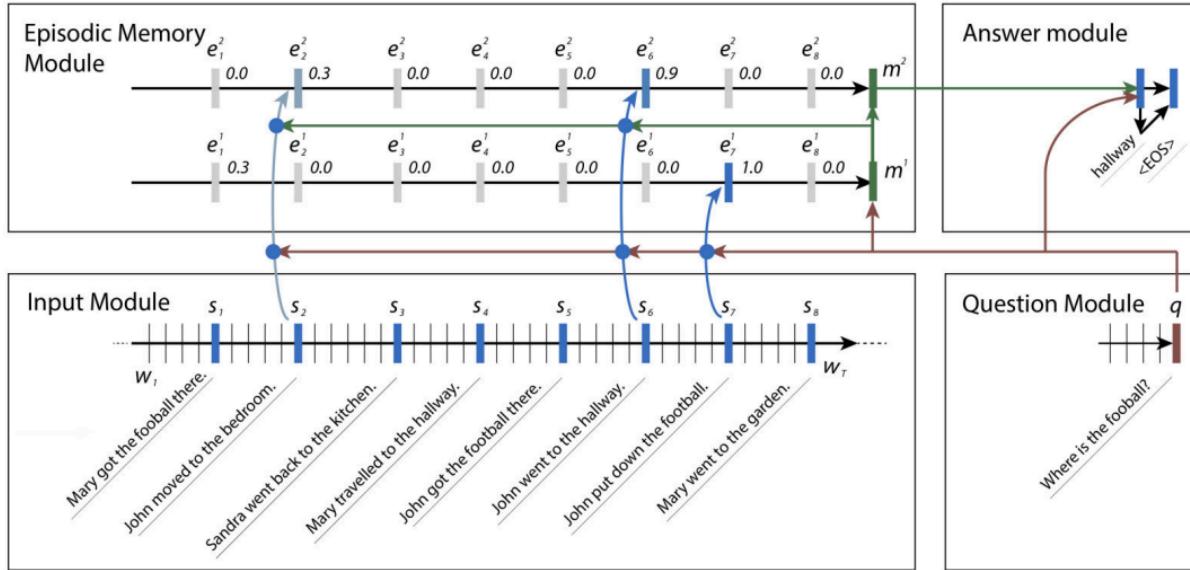
- 谷歌神经翻译 GNMT：Seq2Seq 模型



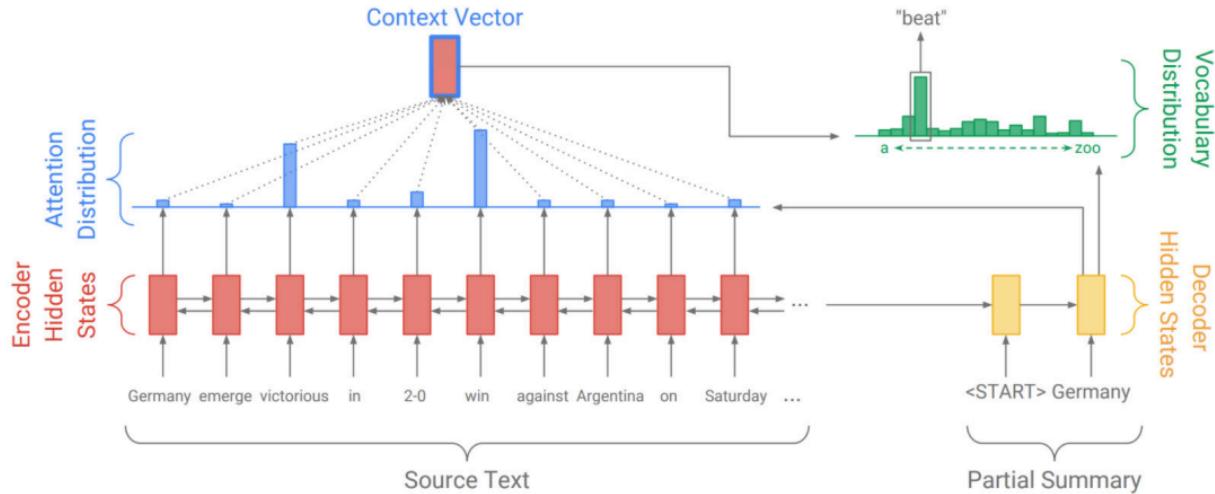
- Google 的 TTS 系统: Tacotron



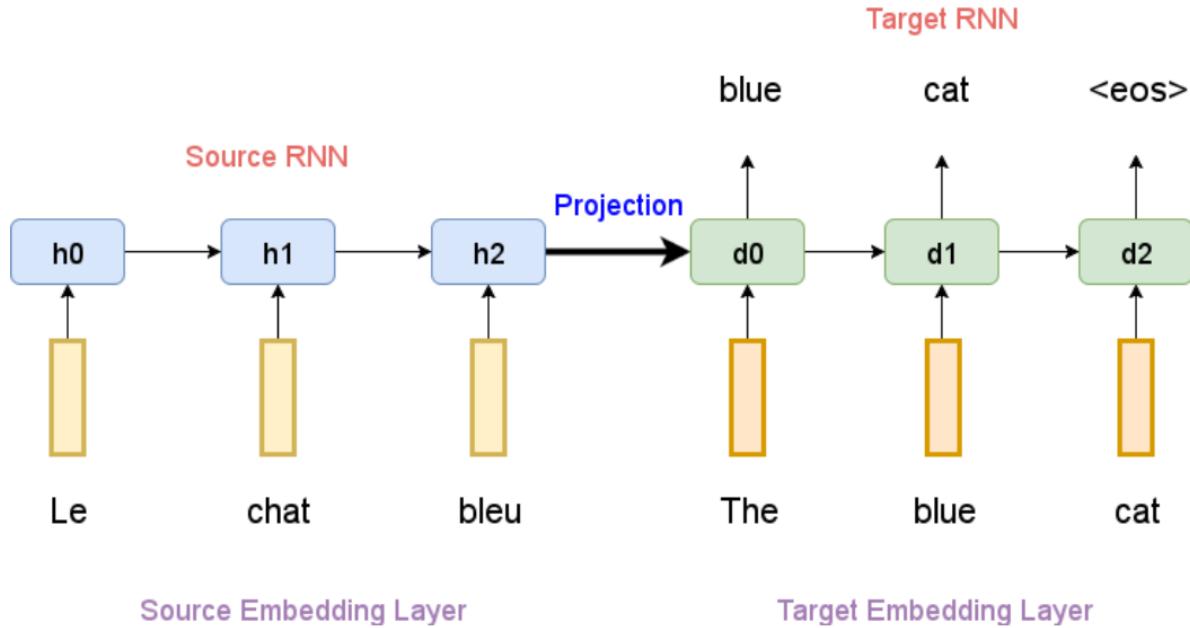
- 问答系统：



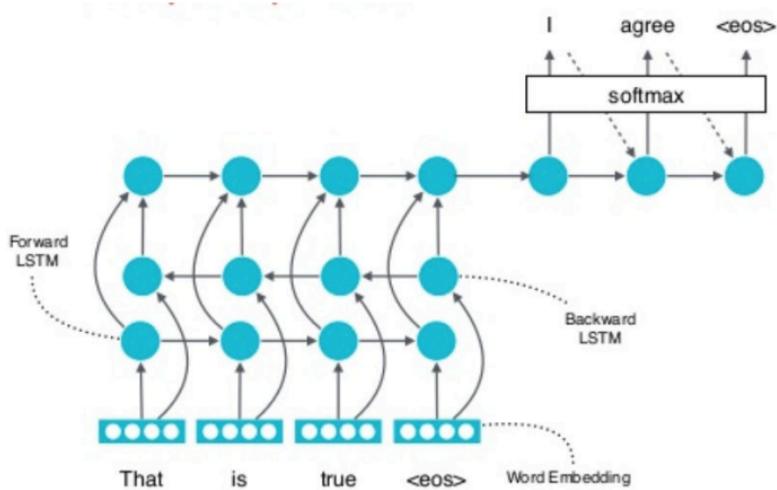
- 概要提取



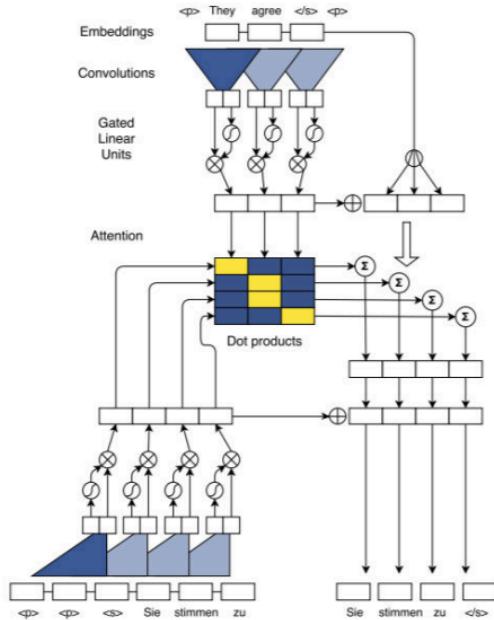
- 最简 Seq2Seq



- Seq2Seq: RNN 强化编解码



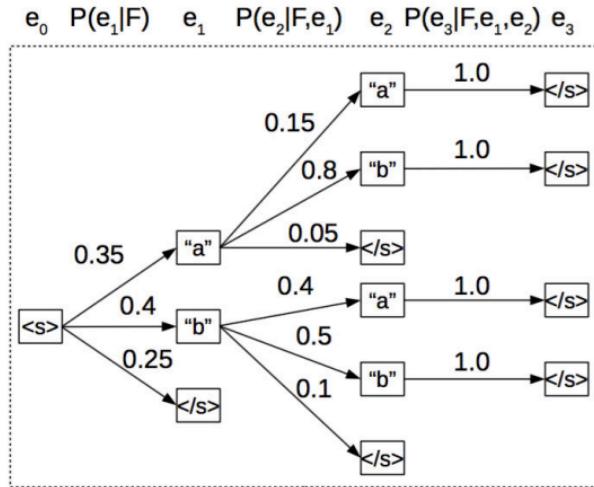
- Seq2Seq: Attention 强化对齐



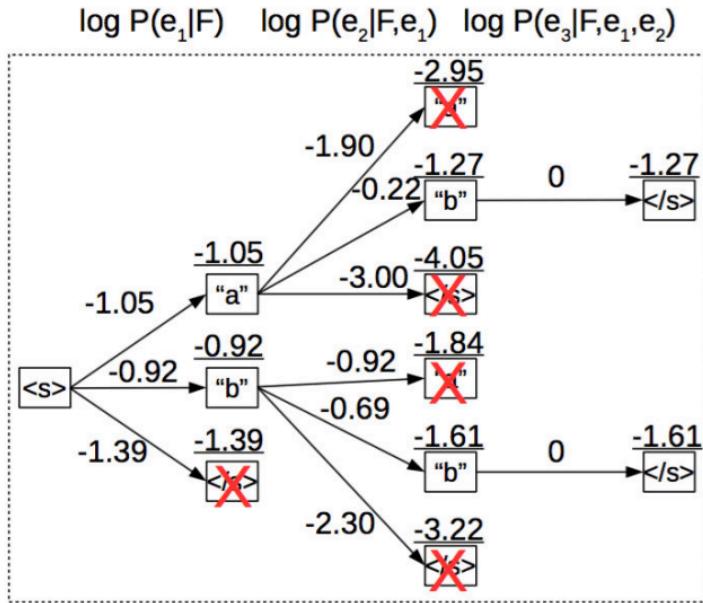
- Seq2Seq 模型小结
  1. Seq2Seq 模型适用范围很广
    - 翻译、问答、概要、语音、图片
  2. Seq2Seq 增强
    - RNN 变型
    - Attention 机制
  3. Seq2Seq 计算量很大

## 4.6 翻译输出

- 贪心查找 Greedy Search 全局最优

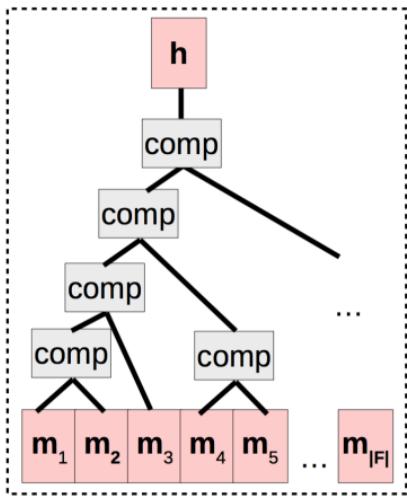


- 束查找 Beam Search 全局最优



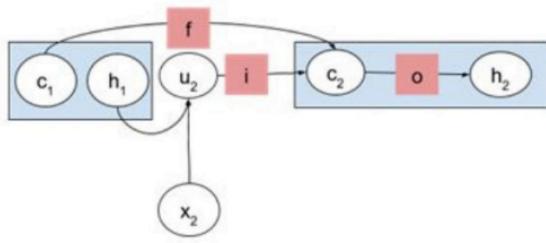
## 4.7 替代和增强 RNN 模型

- 树结构做特征提取

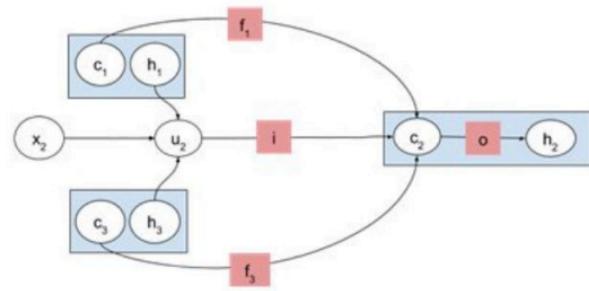


- Tree RNN

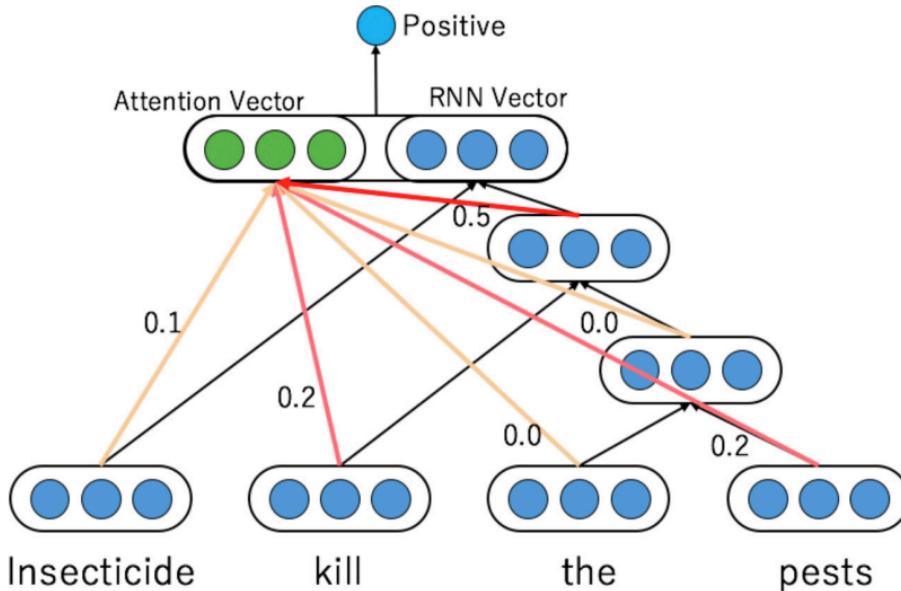
Ordinary Tree-LSTM



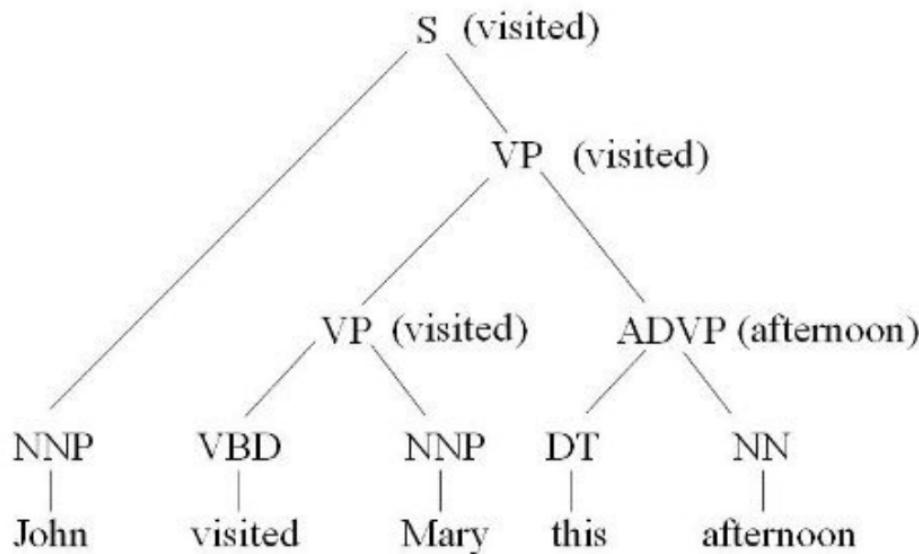
Binary Tree-LSTM



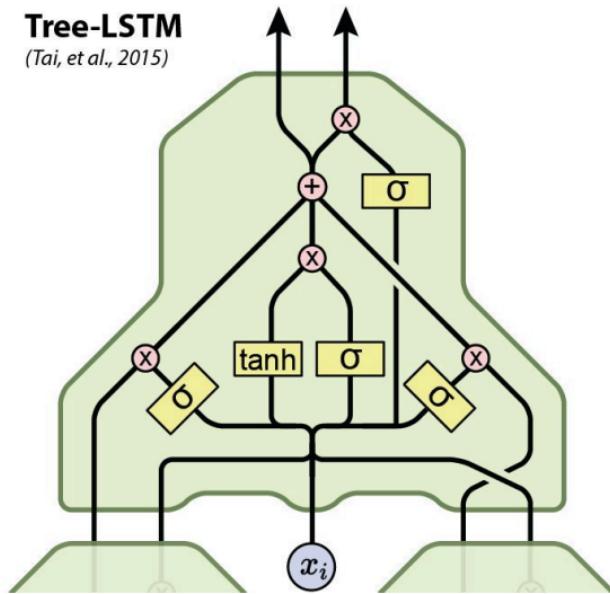
- 二叉树 RNN



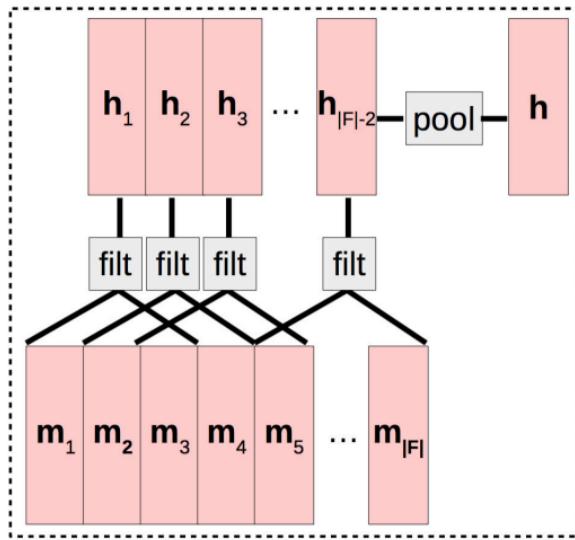
- 语法树



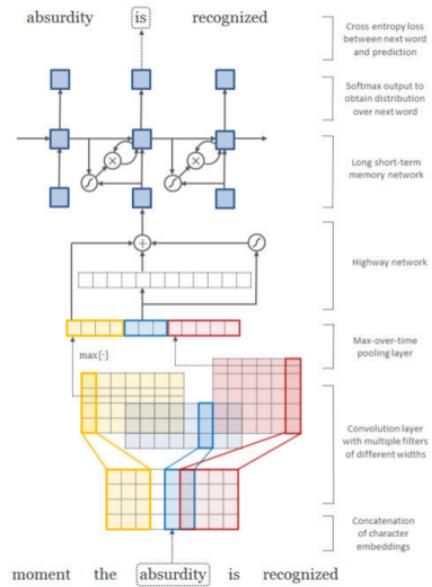
- 树结构化的 LSTM



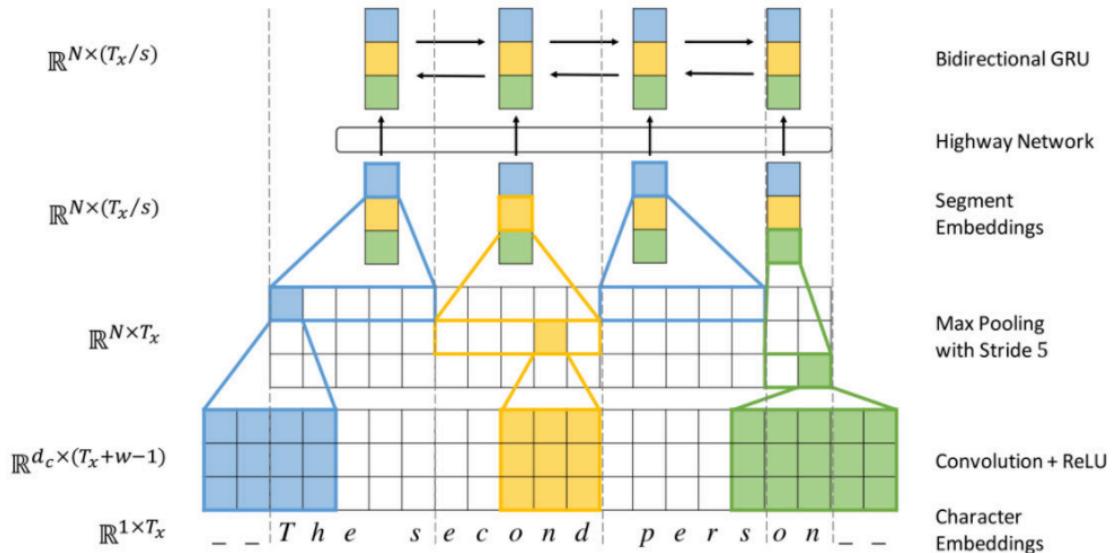
- CNN 做特征提取



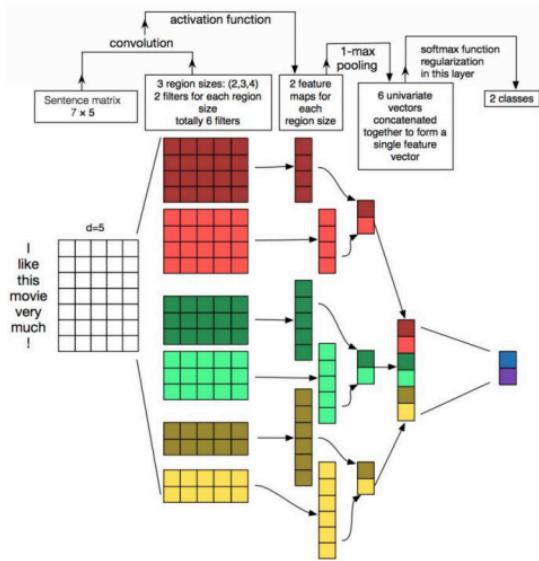
- CNN-RNN 模型: char-cnn-rnn 文本



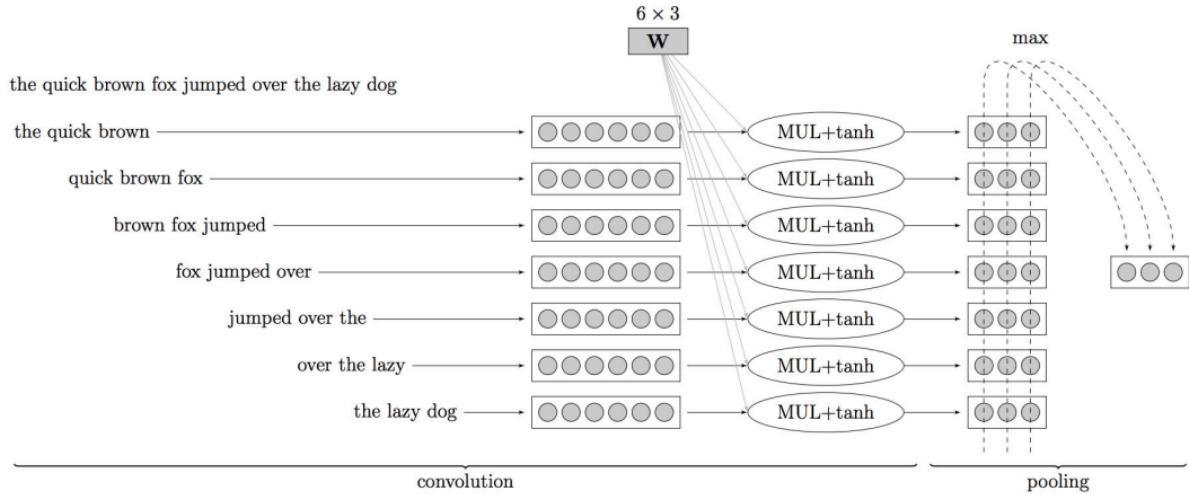
- Char-CNN-RNN 模型



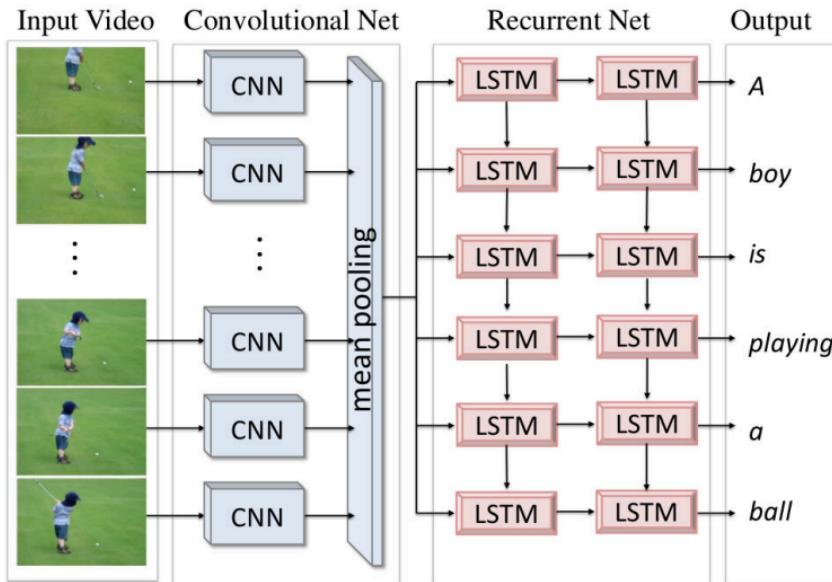
- CNN 文本处理



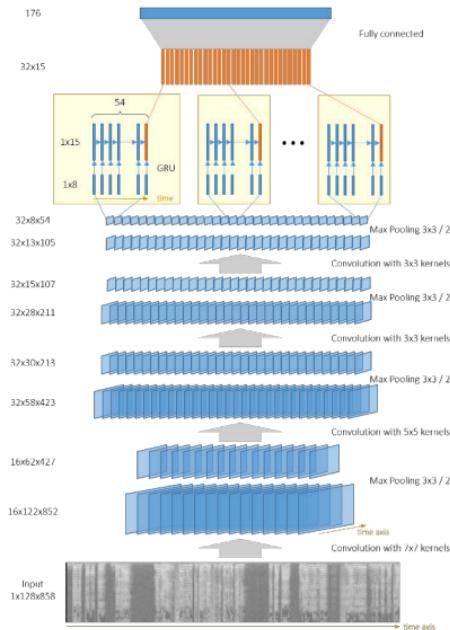
- CNN 文本处理：窗口滑动



- CNN-RNN 模型：图像处理



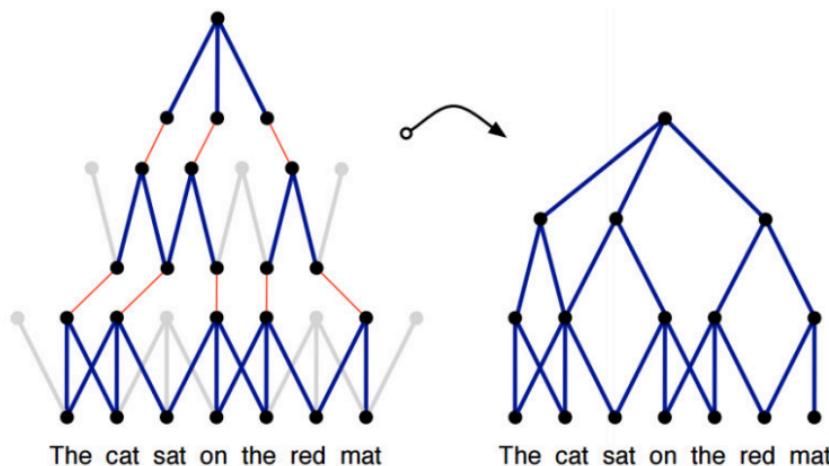
- CNN-RNN 模型：语音处理



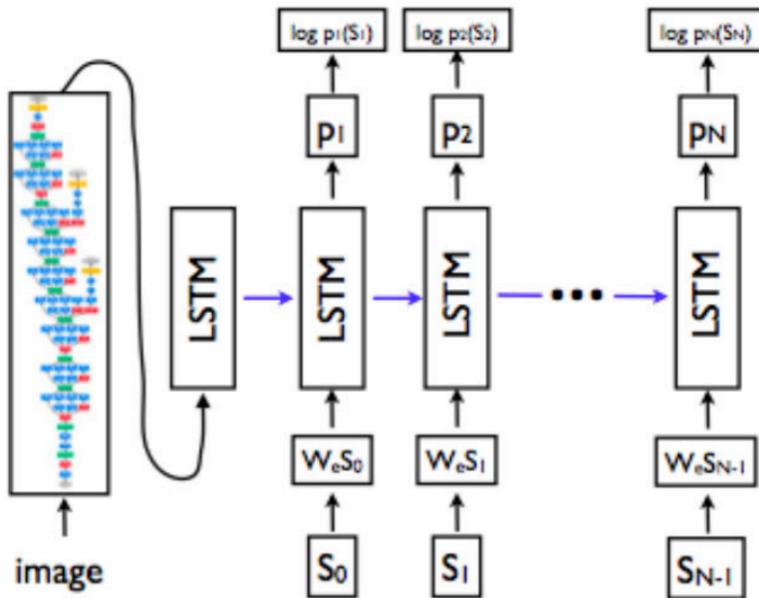
- Tree-RNN 和 CNN-RNN 小结
  1. 通过结构化 RNN 引入结构先验
  2. 通过 CNN 特征提取简化计算
  3. 通过 Tree-RNN 和 CNN-RNN 强化稳定性

## 4.8 其他应用

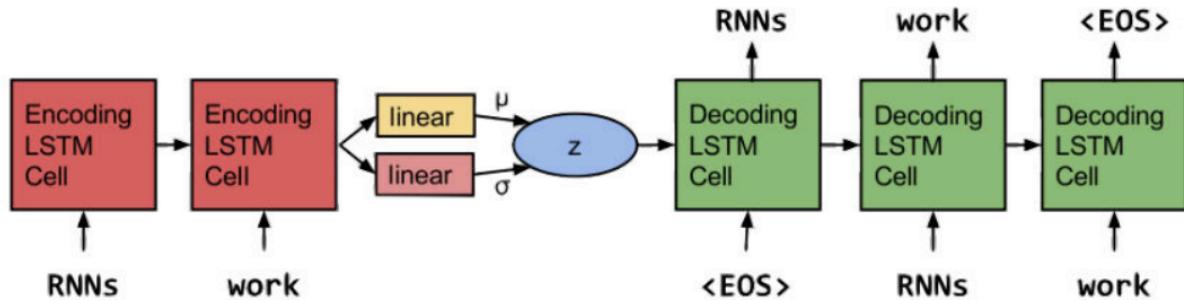
- 句子语言建模：DCNN (dynamic convolutional neural network )



- 图像标题：CNN-LSTM Decoder



- 句子生成：RNN-based VAE



感谢 Stanford, CMU, MIT 等网上公开课程和资料！

# AI2ML

