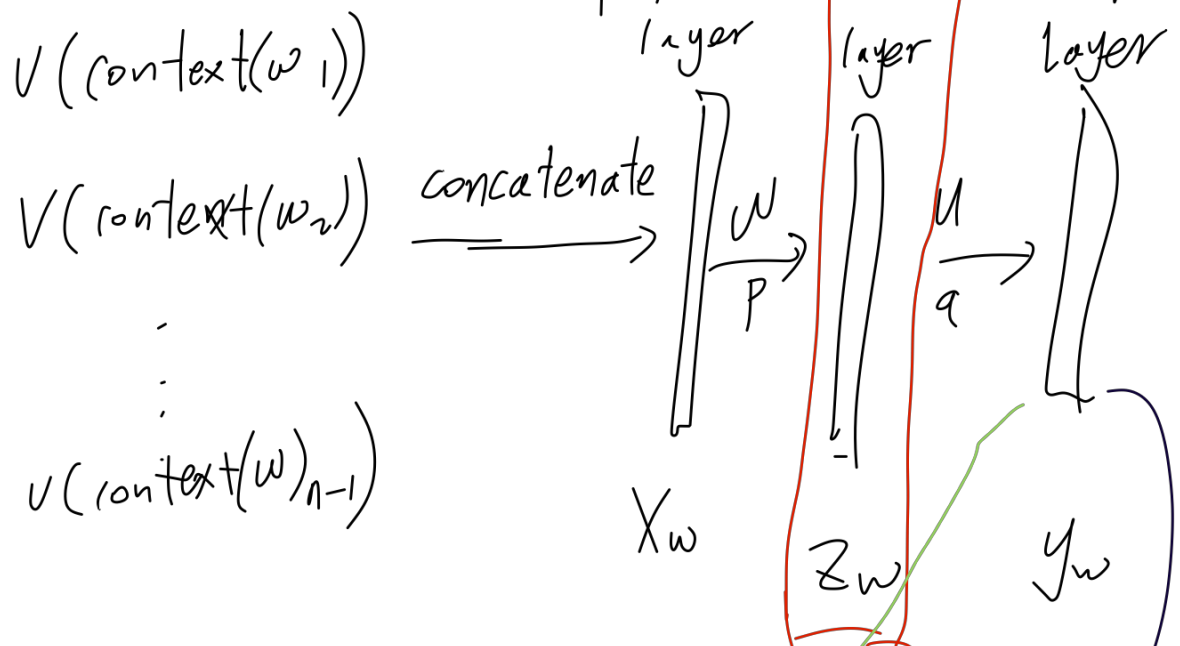


# 基本的神经网络结构



wordvec 去除

hierarchical softmax

对几十万个单词做  
softmax

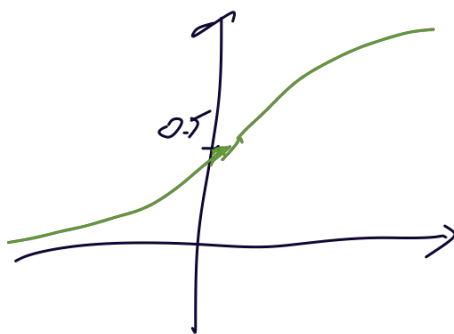
$$s_i = \frac{e^{z_i}}{\sum_j e^{z_j}}$$

得到每个单词的概率

## 2.1 sigmoid

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

$$\sigma'(x) = \sigma(x)(1 - \sigma(x))$$



$$\begin{aligned} \log'(\sigma(x)) &= 1 - \sigma(x) \\ \log'(1 - \sigma(x)) &= -\sigma(x) \end{aligned}$$

逻辑回归做 = 分类, 用了 sigmoid.

$$x = (x_1, x_2, \dots, x_n)^T$$

$$\begin{aligned} h_{\theta}(x) &= (\theta_0 + \theta_1 x_1 + \theta_2 x_2 \dots \theta_n x_n) \\ &= (\theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n) \quad (x_0 = 1) \end{aligned}$$

?

$$h_{\theta}(x) = \sigma(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}}$$

样本  $x$  的取值  $h_{\theta}(x)$

$$y = \begin{cases} 0, & h_{\theta}(x) < 0.5 \\ 1, & h_{\theta}(x) \geq 0.5 \end{cases}$$

$\theta$  是参数, 需更新. 损失函数:

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{cost}(x_i, y_i)$$

$$\text{cost}(x_i, y_i) = \begin{cases} -\log h_{\theta}(x_i) & y_i = 1 \\ -\log (1 - h_{\theta}(x_i)) & y_i = 0 \end{cases}$$

$$P(A|B) = \frac{P(A \cap B)}{P(B)} \quad P(B|A) = \frac{P(A \cap B)}{P(A)}$$

Bayes

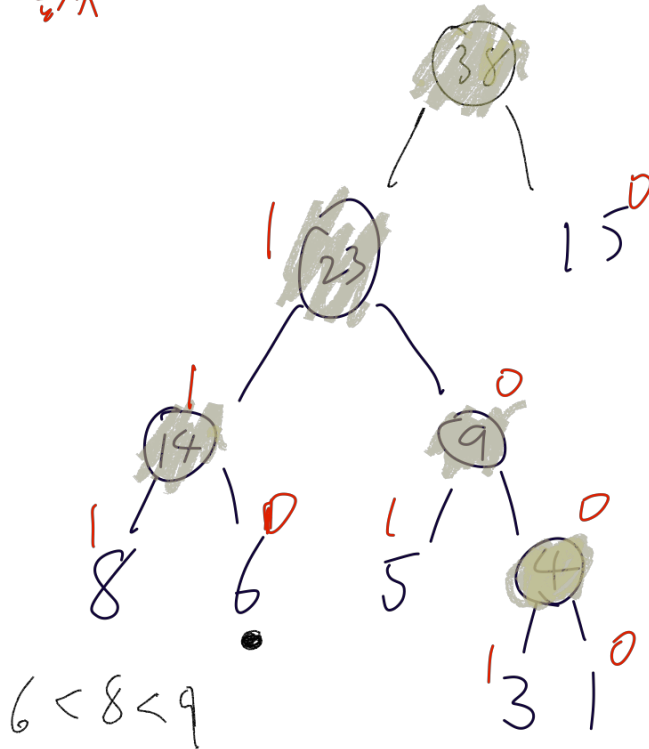
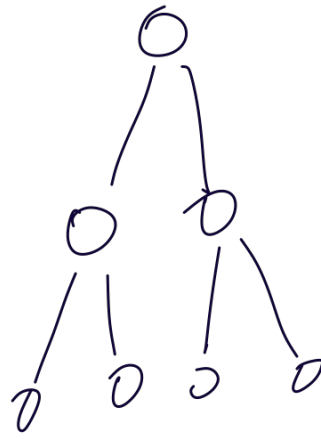
$$\underline{P(A|B)} = \frac{P(A) \cdot \underline{P(B|A)}}{P(B)}$$

两个条件概率

Huffman 树是二叉树。

例子:

15, 8, 6, 5, 3, 1.  
我 喜欢 观看 巴西 足球 世界杯



词频越大的词离根结点越近。

Huffman 编码

是不等长编码，利用霍夫曼树  
设计的二进制前缀码

wordzvel 也用到了霍夫曼编码。

$$P(\text{text}|\text{voice}) = \frac{P(\text{voice}|\text{text}) P(\text{text})}{P(\text{voice})}$$

语言模型

声学模型

$$\begin{aligned} \text{LM: } P(w_1 w_2 \dots w_n) \\ &= P(w_n) \cdot P(w_n | w_1 \dots w_{n-1}) \\ &= P(w_n). \end{aligned}$$

$$P(A|B) = \frac{P(AB)}{P(B)}$$

$$\begin{aligned} P(w_1 w_2 w_3) \\ &= P(w_1) \cdot P(w_2 | w_1) \cdot P(w_3 | w_2). \end{aligned}$$

$$\begin{aligned} P(w_1 w_2 w_3 \dots w_n) \\ &= P(w_1) \cdot P(w_2 | w_1) \cdot P(w_3 | w_2) \dots P(w_n | w_{n-1}) \end{aligned}$$

词典个数为  $N$ . 对于长度为  $T$  的句子,

每个位置有  $N$  种可能,  $\underbrace{N \cdot N \cdot N \dots N}_T = N^T$

每种可能要计算  $T$  个参数, 共需  $T \cdot N^T$

内存开销很大

$n$ -gram, 决策树, 最大熵模型, 最大熵模型,  
条件随机场, 神经网络. 构建语言模型.

$n$ -gram 模型.

$O(N^n)$  词典大小  $N=20000$   
 $n$  是  $n$  gram.  $= 1, 2, 3, 4$ .

$$p(w_1, w_2, w_3, w_4, w_5) \quad 3\text{-gram}$$
$$= p(w_1) \cdot p(w_2/w_1) \cdot p(w_3/w_1, w_2) \cdot p(w_4/w_2, w_3) \\ \cdot p(w_5/w_3, w_4) \quad ? \text{ 对不对.}$$

对问题建模, 构造目标函数, 进行优化.

最大似然

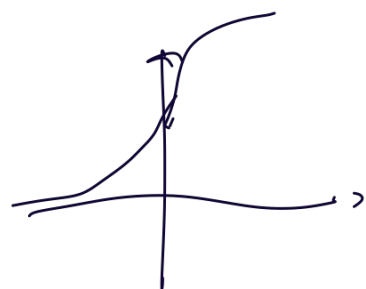
$$\prod_{w \in L} P(w | \text{context}(w))$$

最大对数似然

$$L = \sum_{w \in L} \log P(w | \text{context}(w))$$



tanh 双曲正切函数, 作为隐层的激活函数



词向量

one hot representation  $\rightarrow 0/1$   
distribution representation  
 $\downarrow$   
将概率分配给每个单词

word2vec  $\rightarrow$  CBow  $\rightarrow$  hierarchical softmax  
 $\downarrow$  skip-gram  $\rightarrow$  negative sampling

CBOW

目标函数:

$$L = \sum_{w \in L} \log P(w | \text{context}(w))$$

相对于神经网络的语言模型,  
1' 去掉前隐藏层.

2' 用 Huffman 树替代 softmax.

$$P(w | \text{context}(w)) = \prod_{j=2}^L P(d_j^w | x_w, \theta_{j-1}^w)$$

故  $L$  是关于  $\theta$  和  $w$  的函数.

求  $L$  关于  $\theta$  和  $w$  的偏导. 并更新.

$$\theta_{j-1}^w := \theta_{j-1}^w + \eta \left[ \frac{\partial L(w, j)}{\partial \theta_{j-1}^w} \right]$$

$$V(\tilde{w}) := V(\tilde{w}) + \eta \sum_{j=2}^L \frac{\partial L(w, j)}{\partial w}$$

skip-gram 目标函数

$$L = \sum_{w \in C} \log P(\text{context}(w) | w)$$