

性别	兴趣-钢琴	兴趣-跳舞	兴趣-击剑	兴趣-篮球
1	0.35	0.47	0.90	0.85
1	0.45	0.32	0.67	0.89
1	0.25	0.34	0.88	0.91
0	0.65	0.78	0.12	0.35
0	0.45	0.98	0.38	0.43
?	0.34	0.47	0.60	0.88

$$y_i = w_1 x_1 + w_2 x_2 + w_3 x_3 + w_4 x_4 + b$$

求解 w_1, w_2, w_3, w_4, b

y	x1	x2	x3	x4
[0,1]	0.35	0.47	0.90	0.85
[0,1]	0.45	0.32	0.67	0.89
[0,1]	0.25	0.34	0.88	0.91
[1,0]	0.65	0.78	0.12	0.35
[1,0]	0.45	0.98	0.38	0.43
...

×

w1	w2	b1	b2
0.21	0.64	0.35	0.47
0.15	0.32		
0.65	0.30		
0.25	0.75		

训练参数

$$y_i = w_i x_i + b_i$$

损失函数：

$$loss = \sqrt{(y_{true} - y_{pred})^2}$$

预测值-y ₀	预测值-y ₁	真实值-y ₀	真实值-y ₁	损失 abs
0.94	1.28	0	1	1.22
0.80	1.25	0	1	1.05
0.90	1.21	0	1	1.11
0.42	0.96	1	0	1.54
0.60	1.0	1	0	1.44
...

- 1 构建数据：输入 x 是啥，输出 y 是啥
- 2 构建函数：建立的方程是啥
- 3 损失函数：度量预测 y 与真实 y 的差距
- 4 更新参数：如何求解方程

描述一个人：

姓名	性别	身高	体重	...
小赵	1	176	85	...
小王	1	170	77	...
小刘	0	165	70	...
小陈	0	166	56	...
...

目标： 如何用一组数字描述一个单词或者一个字？

祝大家新年快乐！

字	x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8	...	x_n
祝	1	0	0	0	0	0	0	0	...	0
大	0	1	0	0	0	0	0	0	...	0
家	0	0	1	0	0	0	0	0	...	0
新	0	0	0	1	0	0	0	0	...	0
年	0	0	0	0	1	0	0	0	...	0
快	0	0	0	0	0	1	0	0	...	0
乐	0	0	0	0	0	0	1	0	...	0
！	0	0	0	0	0	0	0	1	...	0

字	x_1	x_2	x_3	...	x_m
祝	0.21	0.12	0.32	...	0.32
大	0.12	0.16	0.53	...	0.45
家	0.43	0.34	0.35	...	0.67
新	0.23	0.79	0.78	...	0.56
年	0.56	0.16	0.81	...	0.94
快	0.59	0.38	0.43	...	0.97
乐	0.78	0.19	0.61	...	0.89
！	0.89	0.16	0.82	...	0.65

One-hot 编码

缺点：存储量大，每个字之间没有任何的关系，只是字的编码是唯一表示的

1 构建数据：输入x是啥，输出 y是啥

祝大家新年快乐！

方法1- Skip-Gram

祝 大 家 新 年
大 家 新 年 快
家 新 年 快 乐
新 年 快 乐 ！

...

x1	y
家	祝
家	大
家	新
家	年
...	...

方法2- CBOW

祝 大 家 新 年
大 家 新 年 快
家 新 年 快 乐
新 年 快 乐 ！

...

x1	x2	x3	x4	y
祝	大	新	年	家
大	家	年	快	新
家	新	快	乐	年
新	年	乐	！	快
...

2 构造函数: $y = w_j (w_i x)$

祝大家新年快乐!

x	x ₁	x ₂	x ₃	x ₄	x ₅	...	x _n
家	0	0	1	0	0	...	0

y	y ₁	y ₂	y ₃	y ₄	y ₅	...	y _n
祝	1	0	0	0	0	...	0

0.12	0.16	...	0.45
0.42	0.21	0.43	0.32
0.43	0.34	...	0.67
...
0.56	0.16	...	0.94

n*100

0.12	0.16	...	0.45
0.78	0.32	0.32	0.67
0.43	0.34	...	0.67
...
0.56	0.16	...	0.94

训练参数

100*n

3 损失函数:

$$loss = F(y_{true}, y_{pred})$$

y	y ₁	y ₂	y ₃	y ₄	y ₅	...	y _n
y _{true}	1	0	0	0	0	...	0
y _{pred}	0.14	0.02	0.01	0.23	0.08	...	0.02

缺点: 预测是2万+的分类

4 优化：负采样- negative sampling

祝大家新年快乐！

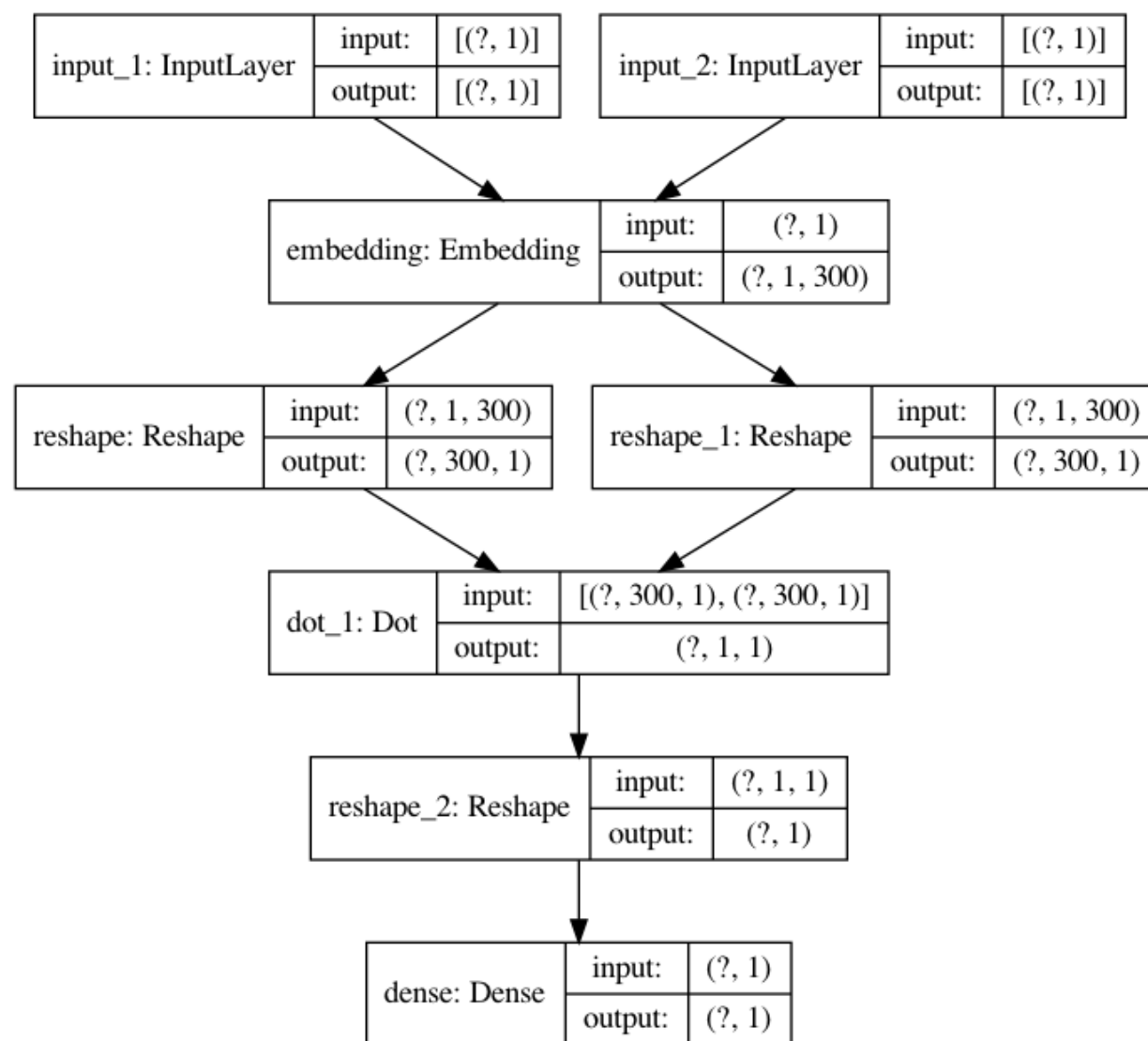
方法1- Skip-Gram

祝大家新年
大家新年快
家新年快乐
新年快乐！
...

x1	y
家	祝
家	大
家	新
家	年
...	...

x1	x2	y
家	祝	1
家	大	1
家	新	1
家	年	1
家	?	0
家	?	0
家	?	0
...

4 优化: 负采样- negative sampling



4 优化：哈夫曼树-Huffman Tree

- 本质是把 N 分类问题变成 $\log(N)$ 次二分类

祝大家新年快乐！

20 38 8 6 16 3 4 25

频率表 快:3 乐:4 新:6 家:8 年:16 祝:20 !:25 大:38

频率表 新:6 快乐:7 家:8 年:16 祝:20 !:25 大:38

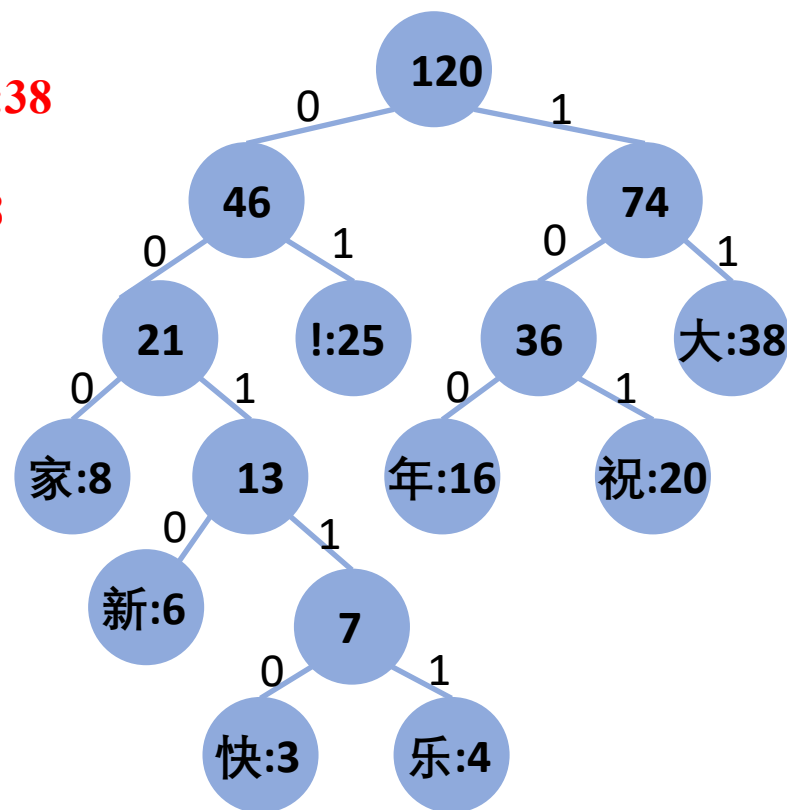
频率表 家:8 快乐新:13 年:16 祝:20 !:25

频率表 年:16 祝:20 快乐新家:21 !:25 大:38

频率表 快乐新家:21 !:25 年祝:36 大:38

频率表 快乐新家!:46 年祝大:74

频率表 快乐新家!年祝大:120



4 优化：哈夫曼树-Huffman Tree

祝:101

大:11

家:000

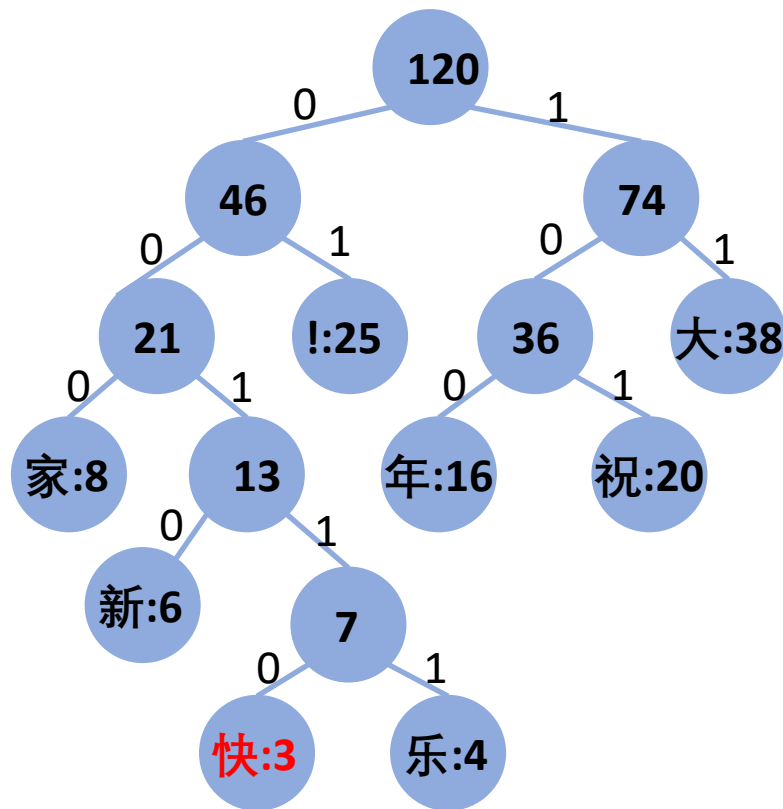
新:0010

年:100

快:00110

乐:00111

! :01



预测的路径概率

$$P = \left(1 - \frac{1}{1 + e^{-x_w \theta_1}}\right) \left(1 - \frac{1}{1 + e^{-x_w \theta_2}}\right) \left(\frac{1}{1 + e^{-x_w \theta_3}}\right) \left(\frac{1}{1 + e^{-x_w \theta_4}}\right) \left(1 - \frac{1}{1 + e^{-x_w \theta_5}}\right)$$