

Word Embedding 是什么?

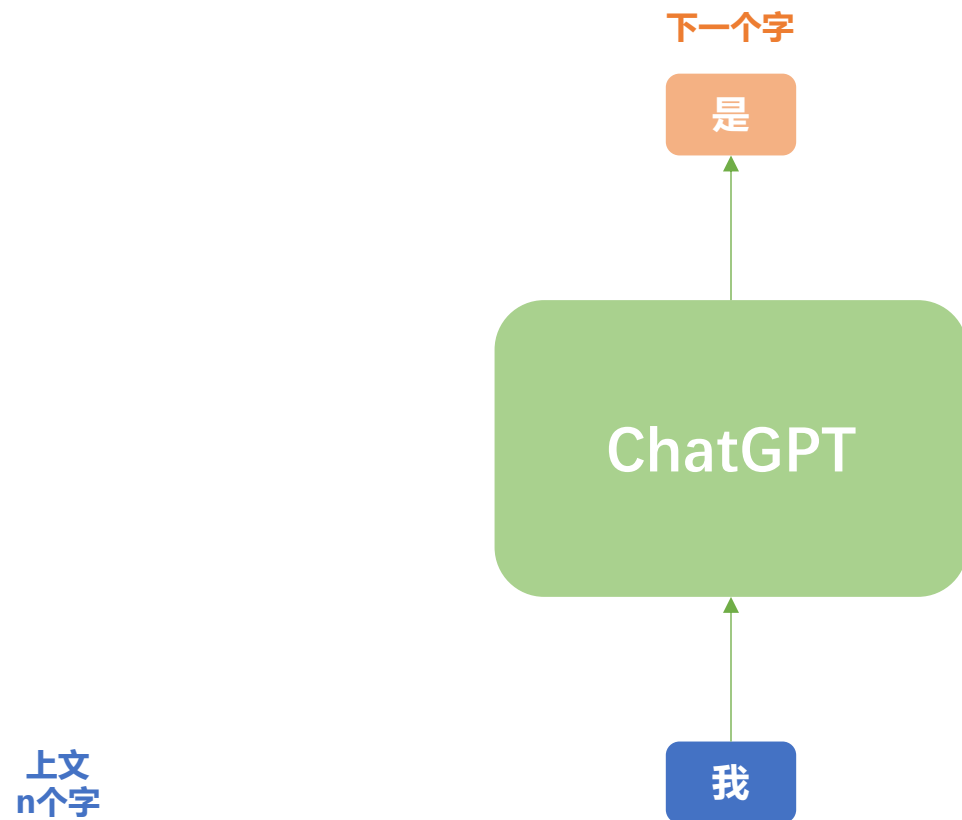
史轩宇

2024.02.24

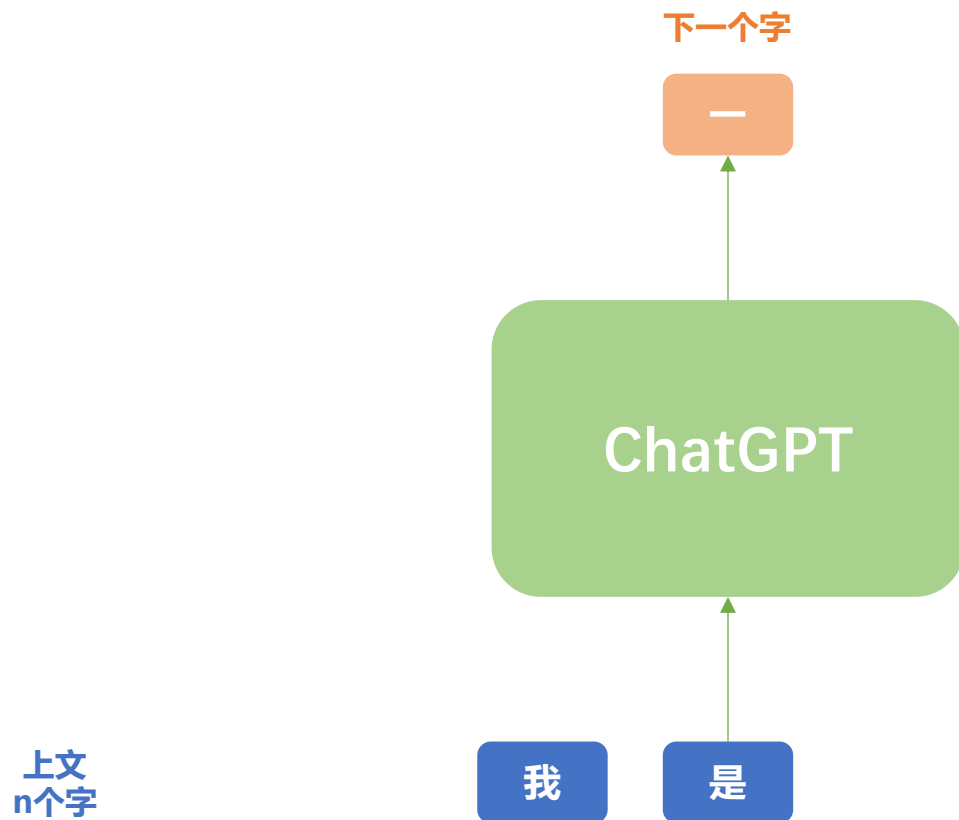


vicshi94

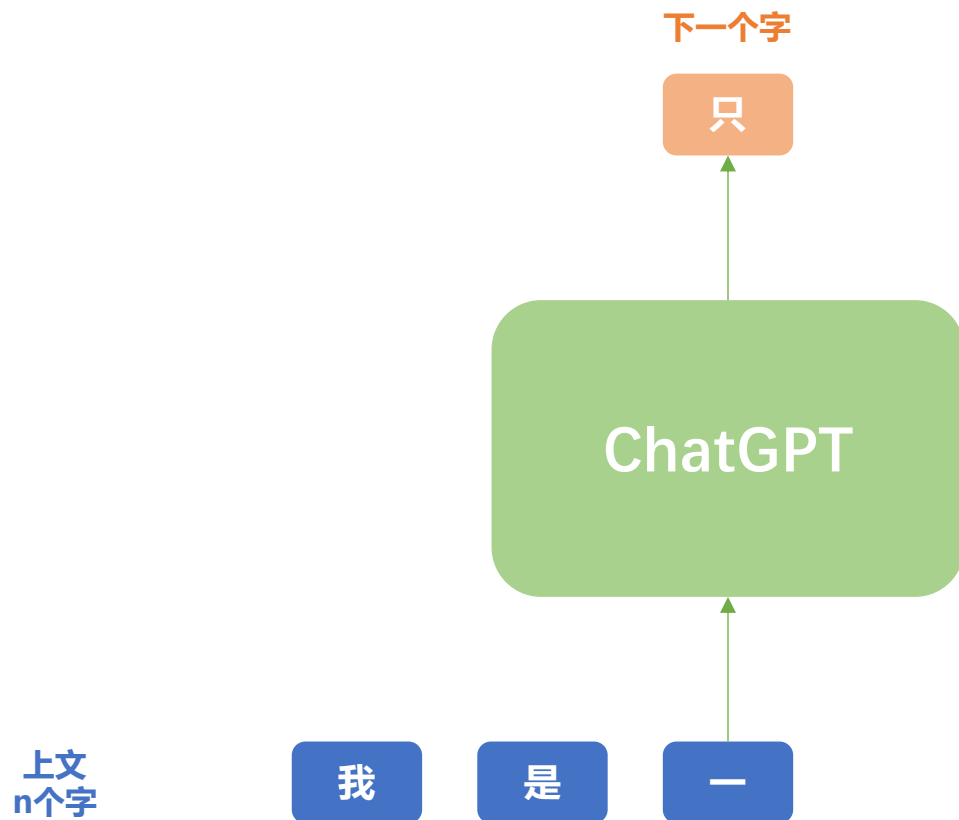
ChatGPT 的基础原理：单字接龙



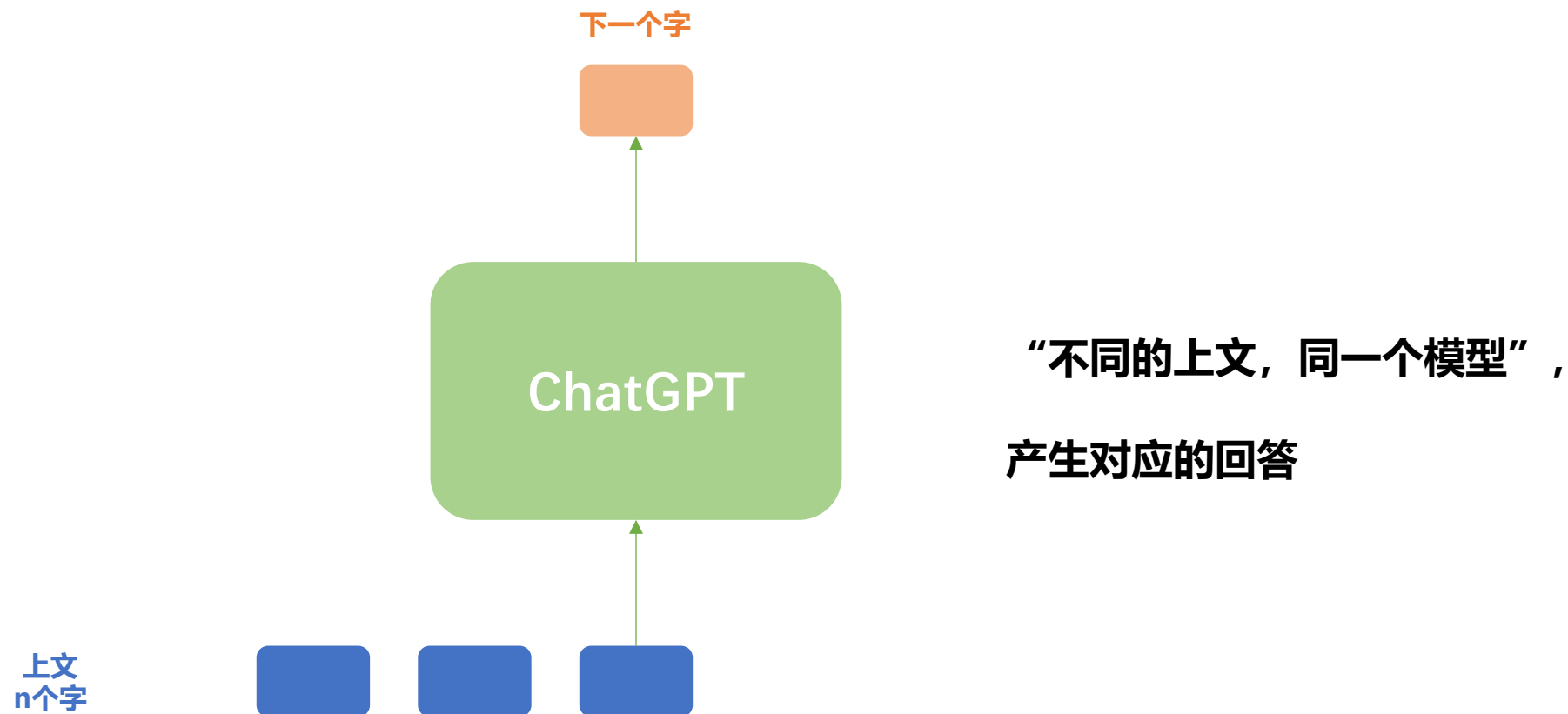
ChatGPT 的基础原理：单字接龙



ChatGPT 的基础原理：单字接龙

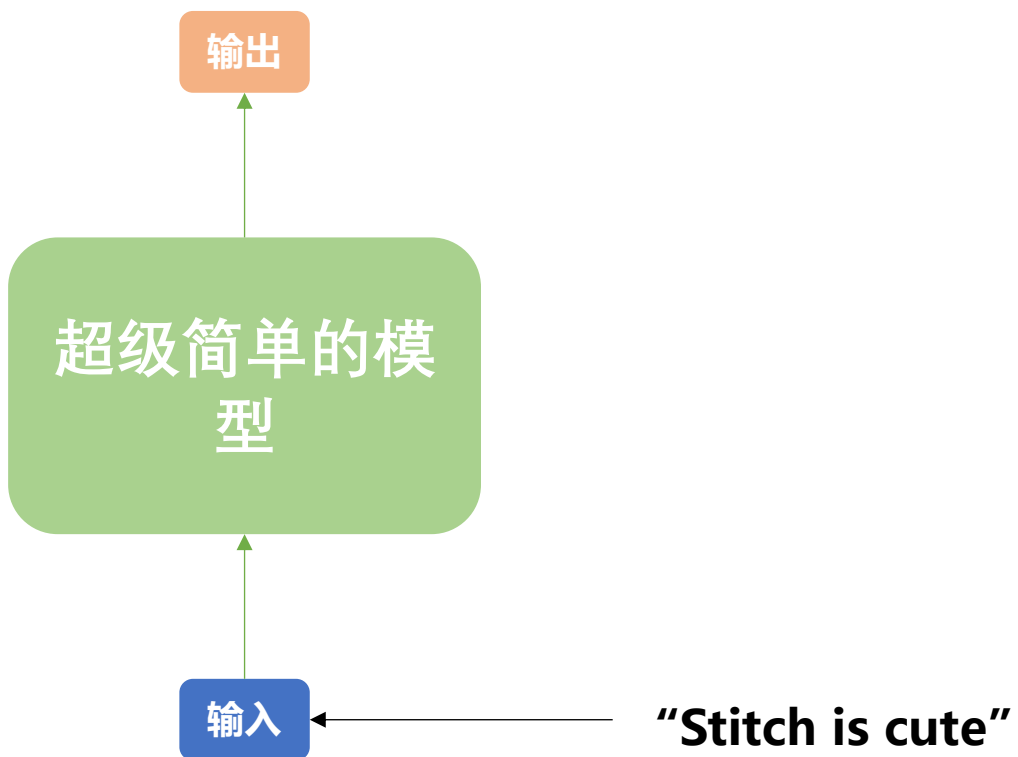


ChatGPT 的基础原理：单字接龙



Text as data?

现在，我们有一个简单模型还有一句话：



怎么把这句话输入进电脑呢？



One-hot Encoding

word	encoding
Stitch	1
is	2
cute	3

An ordinal relationship in numbers!

我们的模型会误以为 “cute” 大于 “is”

	Stitch	is	cute
Stitch	1	0	0
is	0	1	0
cute	0	0	1



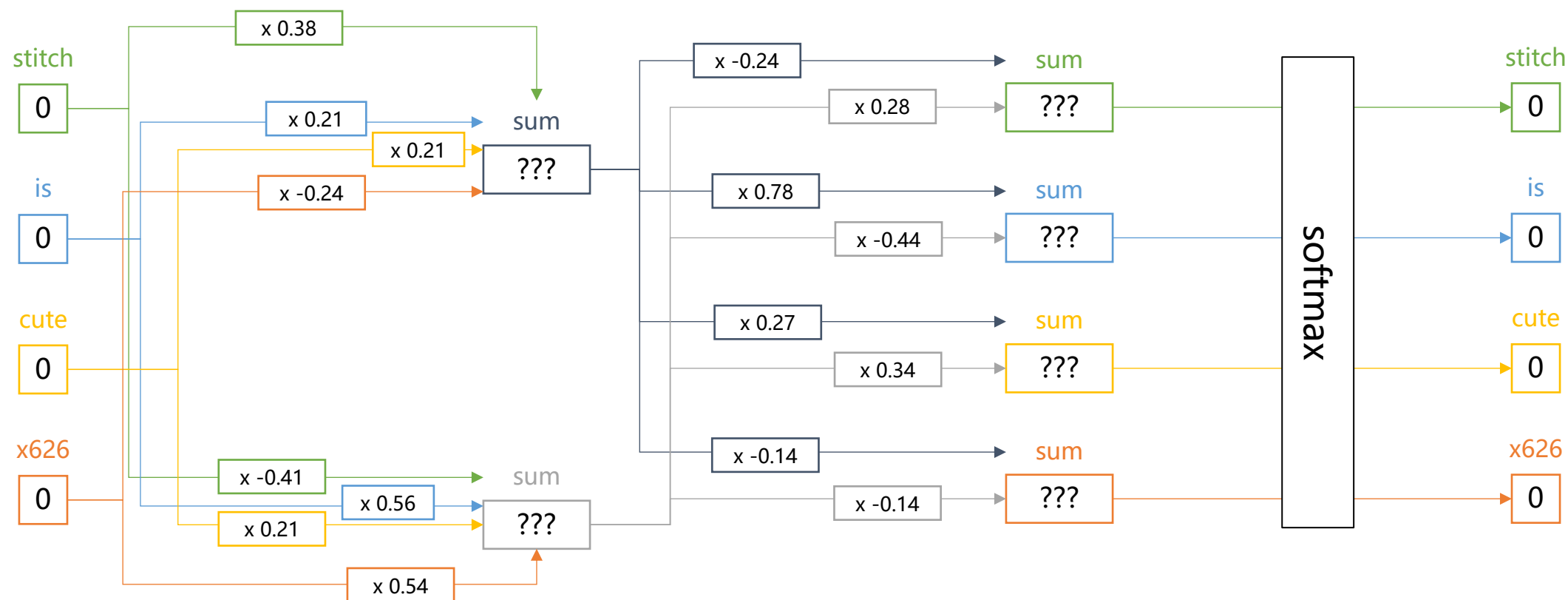
word	one-hot encoding
Stitch	(1,0,0)
is	(0,1,0)
cute	(0,0,1)

这样就没有大小关系了，而且每个单词都被映射为向量空间中距离原点 (0,0,0) 的点，且距离原点的欧式距离都是单位距离1。

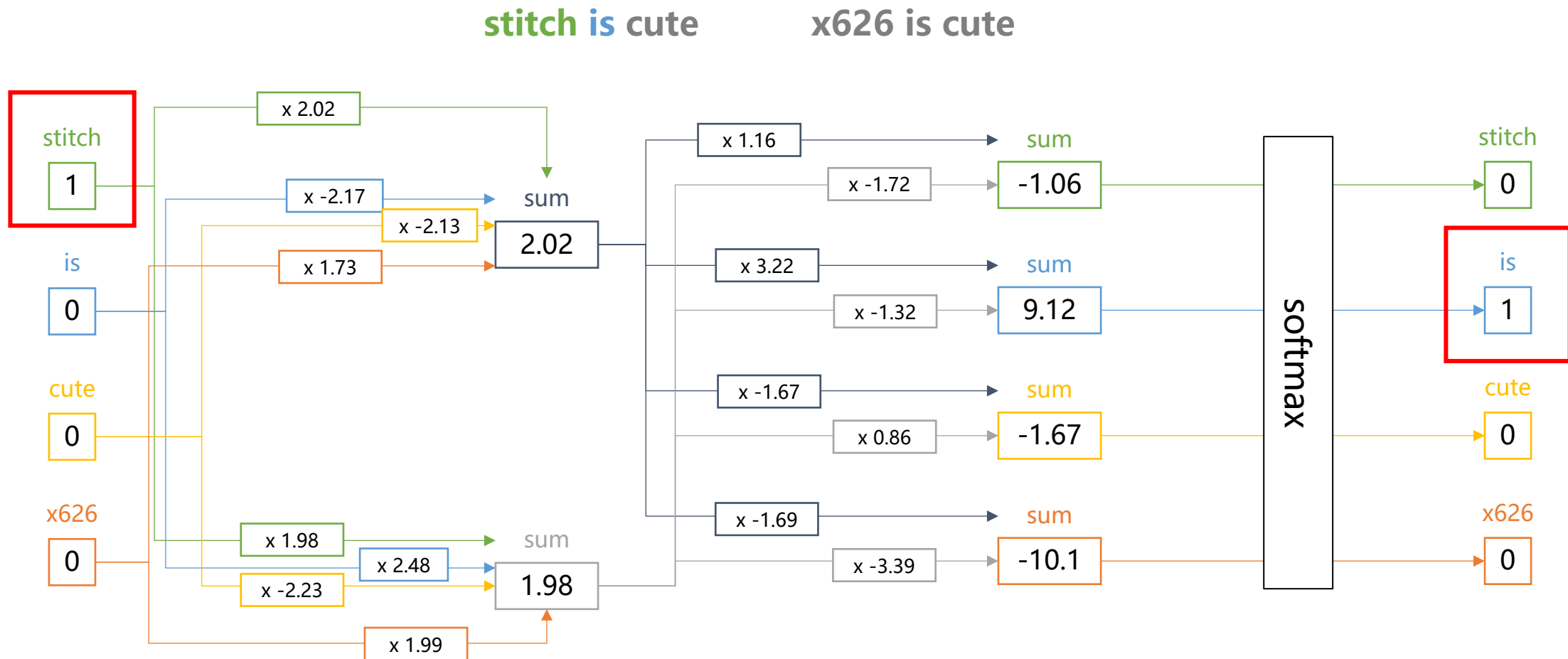
完美！开整！

开始训练超级简单神经网络模型！

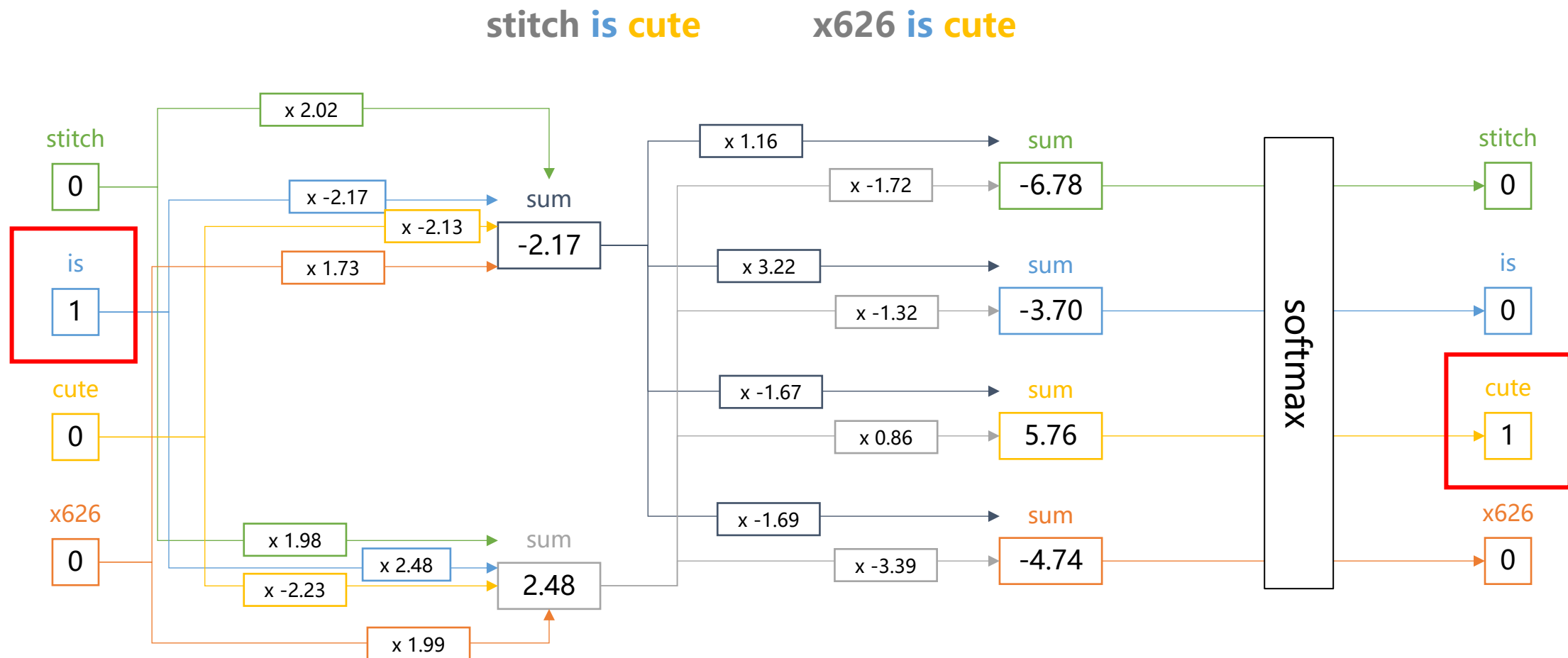
我们的training data是两句话：“stitch is cute”，“x626 is cute”



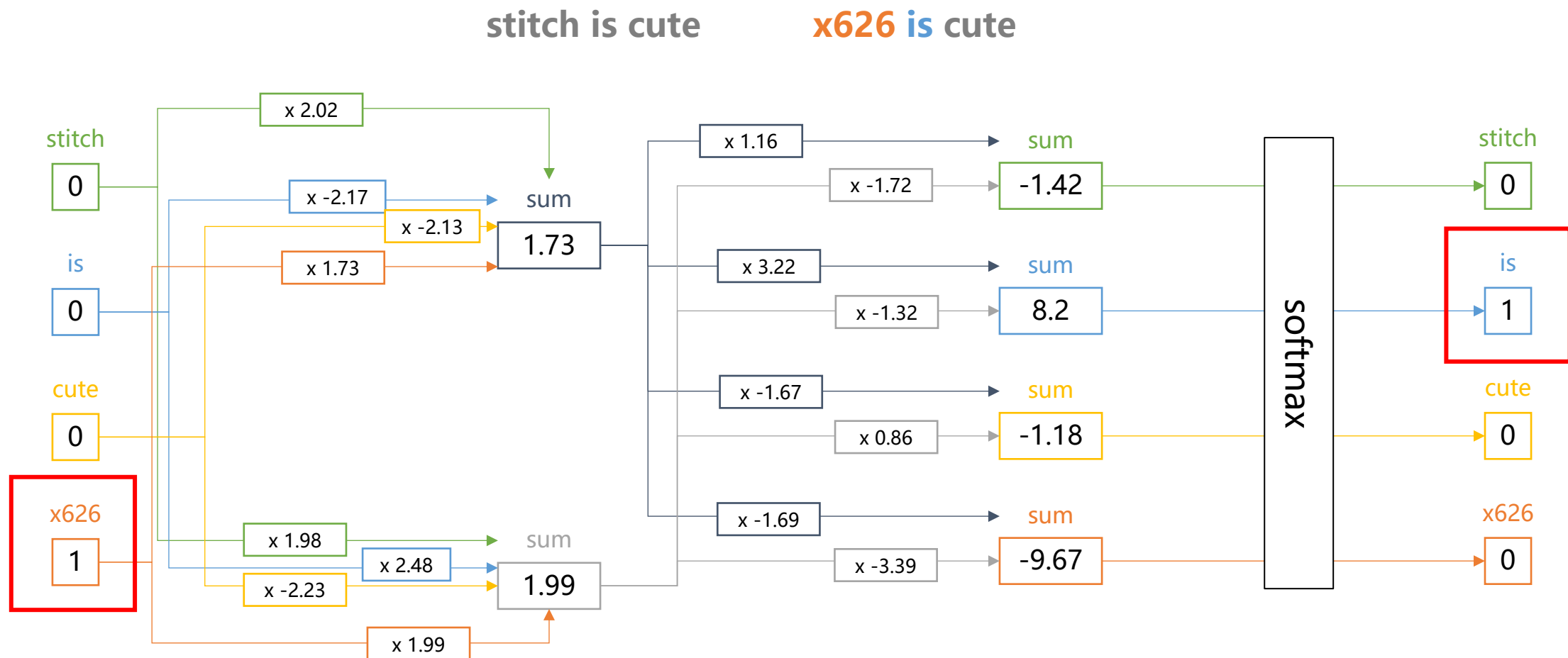
超简单模经过超多轮训练后



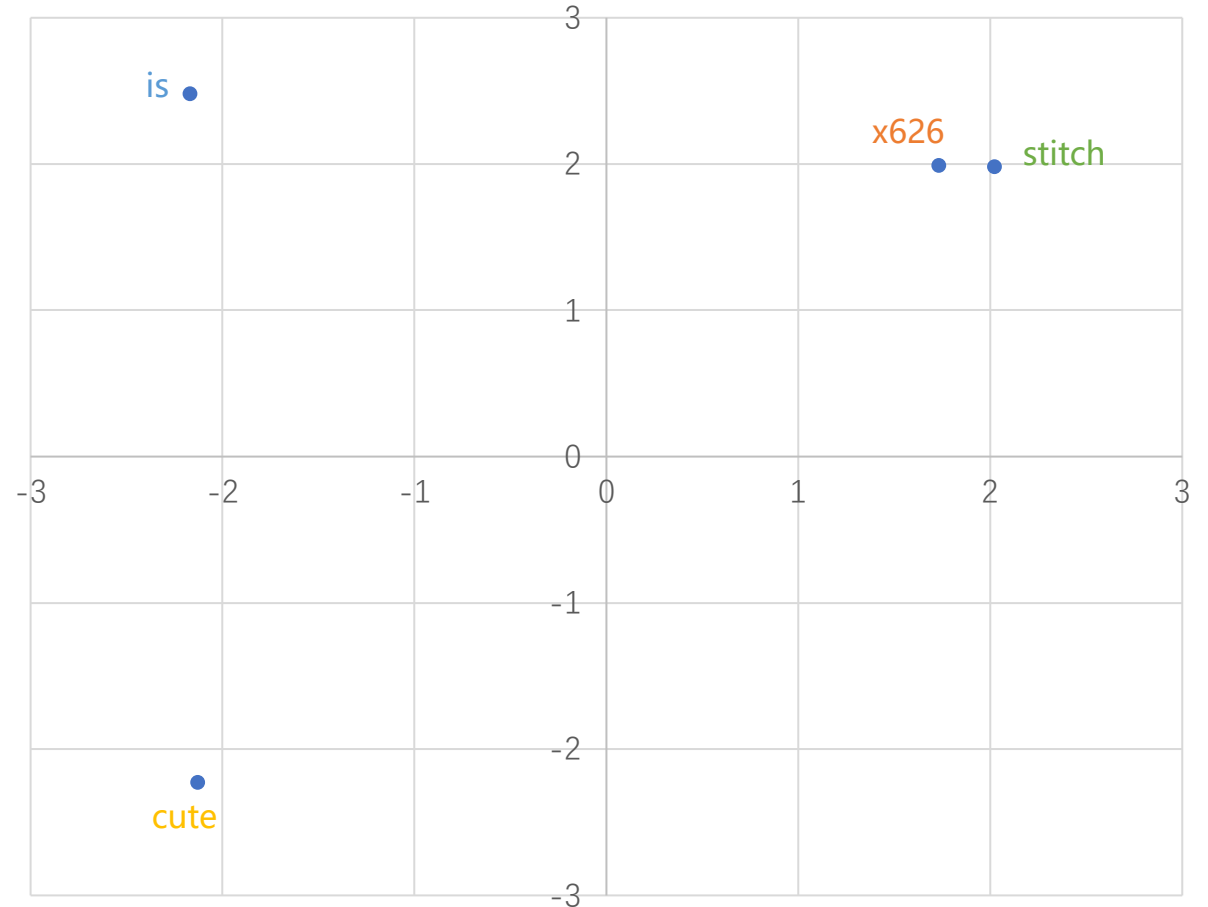
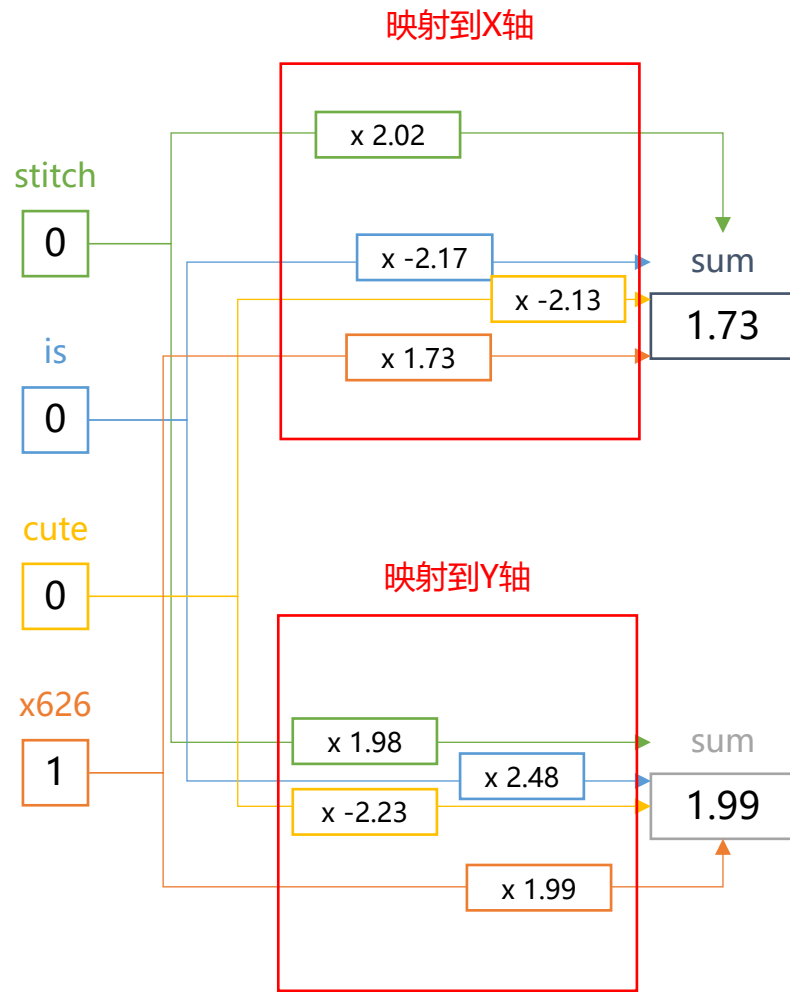
超简单模经过超多轮训练后



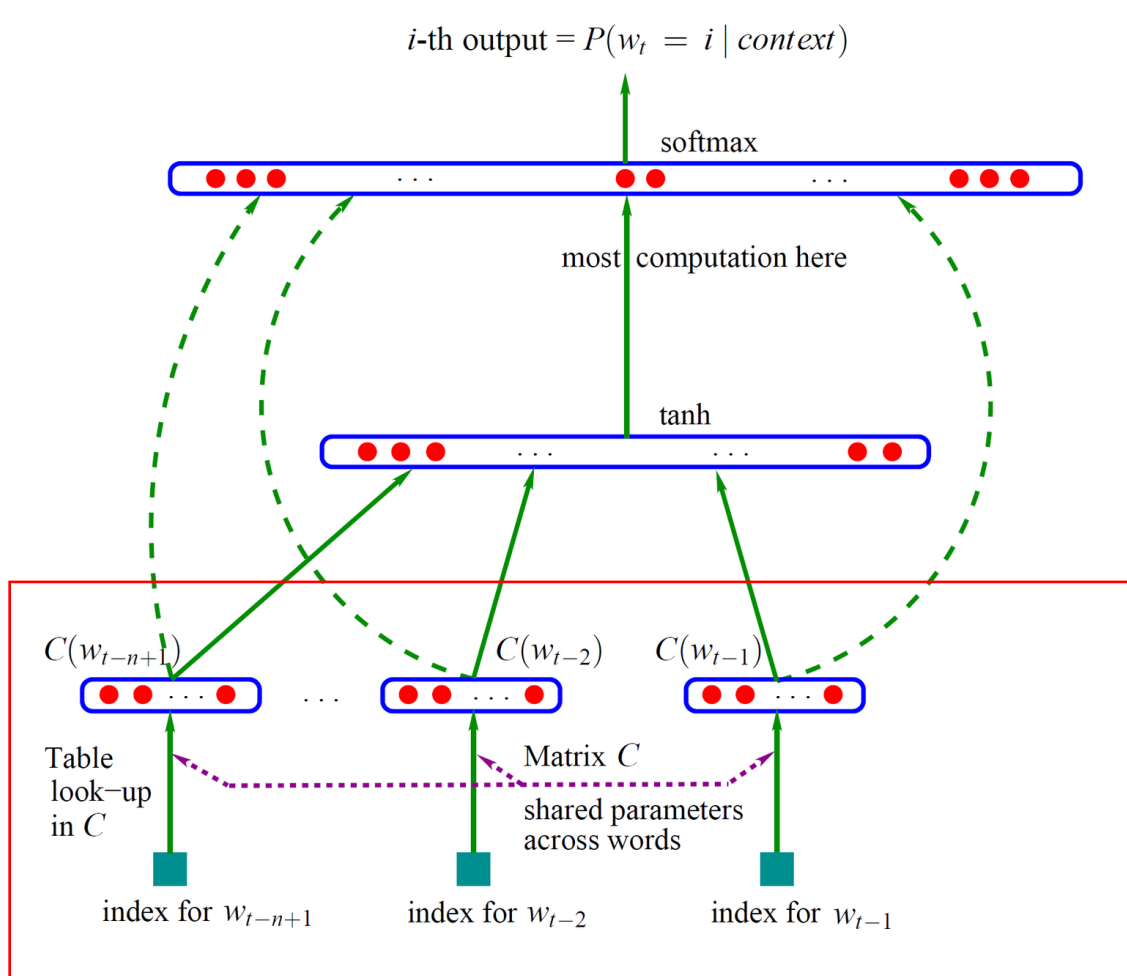
超简单模经过超多轮训练后



超简单模经过超多轮训练后



Vectorization



$$w \times C = c$$

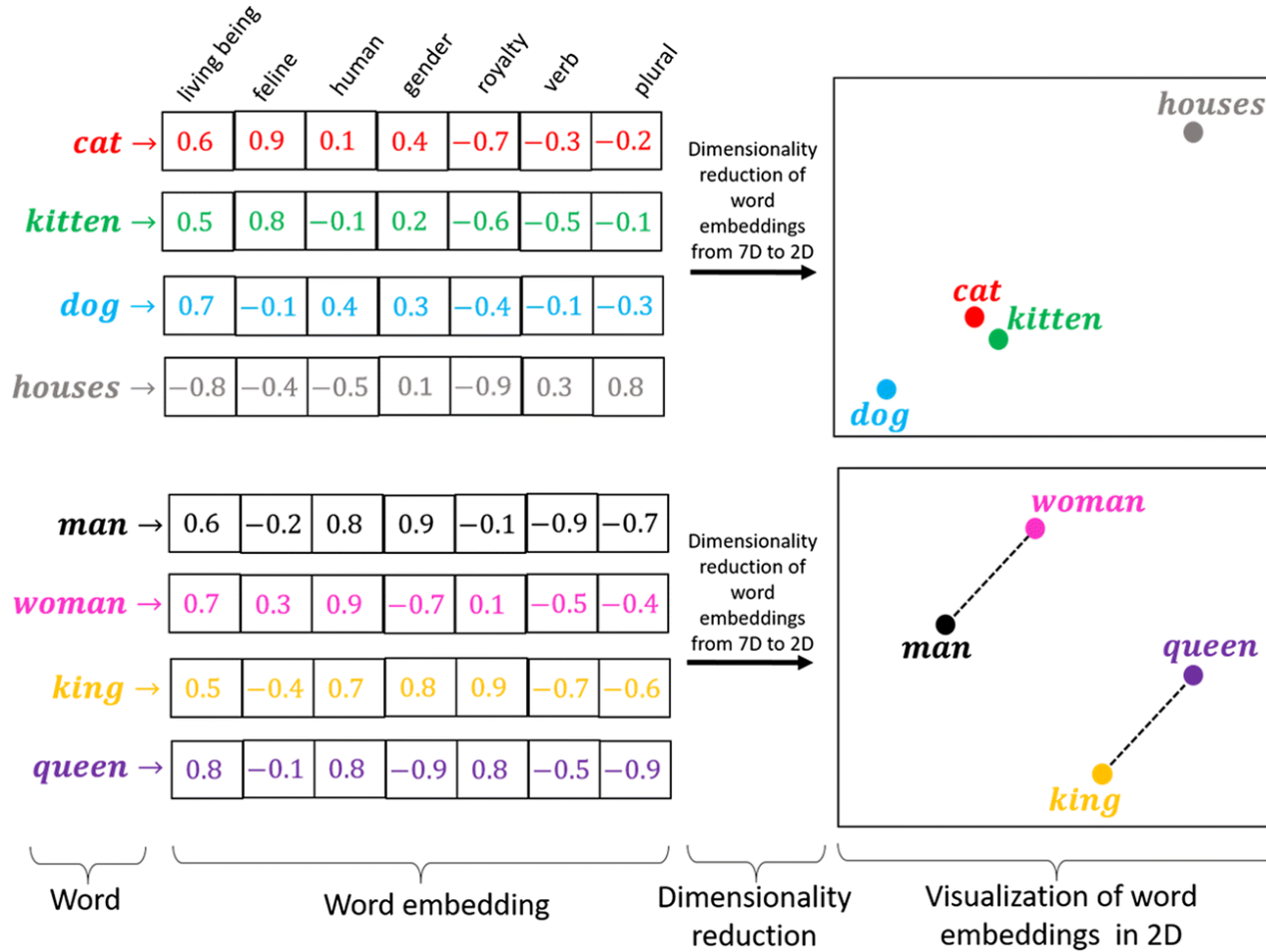
One-hot vector $\begin{bmatrix} 0 & 0 & 0 & 1 & 0 \end{bmatrix} \times$ Embedding Weight Matrix $\begin{bmatrix} 8 & 2 & 1 & 9 \\ 6 & 5 & 4 & 0 \\ 7 & 1 & 6 & 2 \\ 1 & 3 & 5 & 8 \\ 0 & 4 & 9 & 1 \end{bmatrix} =$ Hidden layer output $\begin{bmatrix} 1 & 3 & 5 & 8 \end{bmatrix}$

独热编码主要缺点:

1. 维度太多
2. 太稀疏 ($n^2 - n$ 个 0)

通过word embedding转换为词向量后明显改善了这两点, 且方便了下游任务

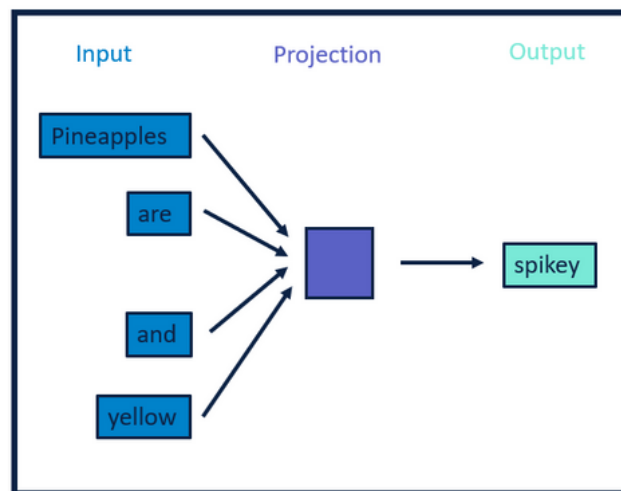
Word Embedding Application



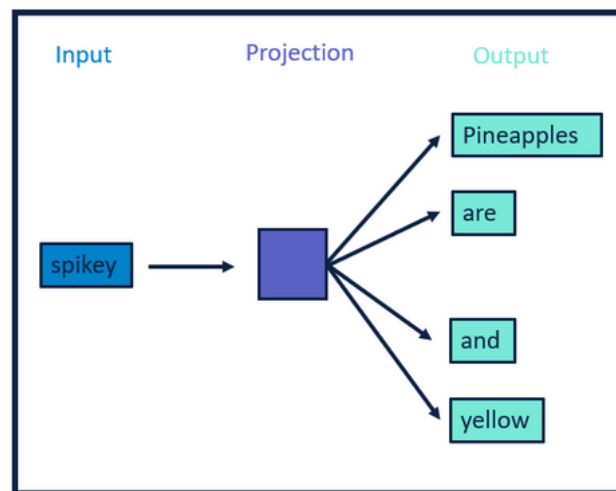
and TM、Supervised ML.....

为了强化对于矩阵C的训练

Word2Vec



CBOW



Skip-gram

GloVe

简单来说，更强调全局的共现频率，而非仅捕捉相邻的几个词

Tutorial: text2vec in R