

简介: word2vec 是主流词嵌入方法, 它不再统计计算词  $w$  在邻域中出现  $s$  的频率, 而是训练一个分类器去预测  $w$  邻域中出现  $s$  的概率, 最终根据分类器的权重构建词向量.

## Skip-gram 模型

任务: 给定目标词, 预测它的上下文

原理: 词相似  $\rightarrow$  上下文相似

不直接算共现, 而是做一个分类器区分 "是否是真的上下文"

最终分类器权重 = 词向量

算法: 1. 定窗口

设窗口大小  $L$  (如  $\pm 2$ ), 取目标词周围的上下文

... lemon, a tablespoon of apricot jam a pinch ...  
                  c1      c2 target c3   c4

2. 构造正负样本

正样本: 选邻域一个词作正样本 (真正的上下文)

负样本: 在词表中随机选一个词作负样本 (噪声)

3. 训练二分类器

给定一个词对  $(t, c)$   $t = \text{target}$   $c = \text{context}$

训练一个分类器, 返回  $c$  是  $t$  上下文的概率

$c$  是上下文:  $P(+|w, c)$

$c$  不是上下文:  $P(-|w, c) = 1 - P(+|w, c)$

注:  $w$  是  $t$  的词向量,  $c$  是  $c$  的词向量

在训练开始前, 这两个向量的值是随便赋的

如何计算  $P(+|w, c)$ ? 用向量内积度量相似, 即  $\text{Similarity}(w, c)$ , 而  $\text{Similarity}(w, c) \propto w \cdot c$ , 直接简化为

$w \cdot c$  (实际上 Word2vec 会对向量作归一化处理, 所有向量模长都为 1) 然后再用 sigmoid 函数转化为概率, 即

$$P(+|w, c) = \sigma(c \cdot w) = \frac{1}{1 + \exp(-c \cdot w)}$$

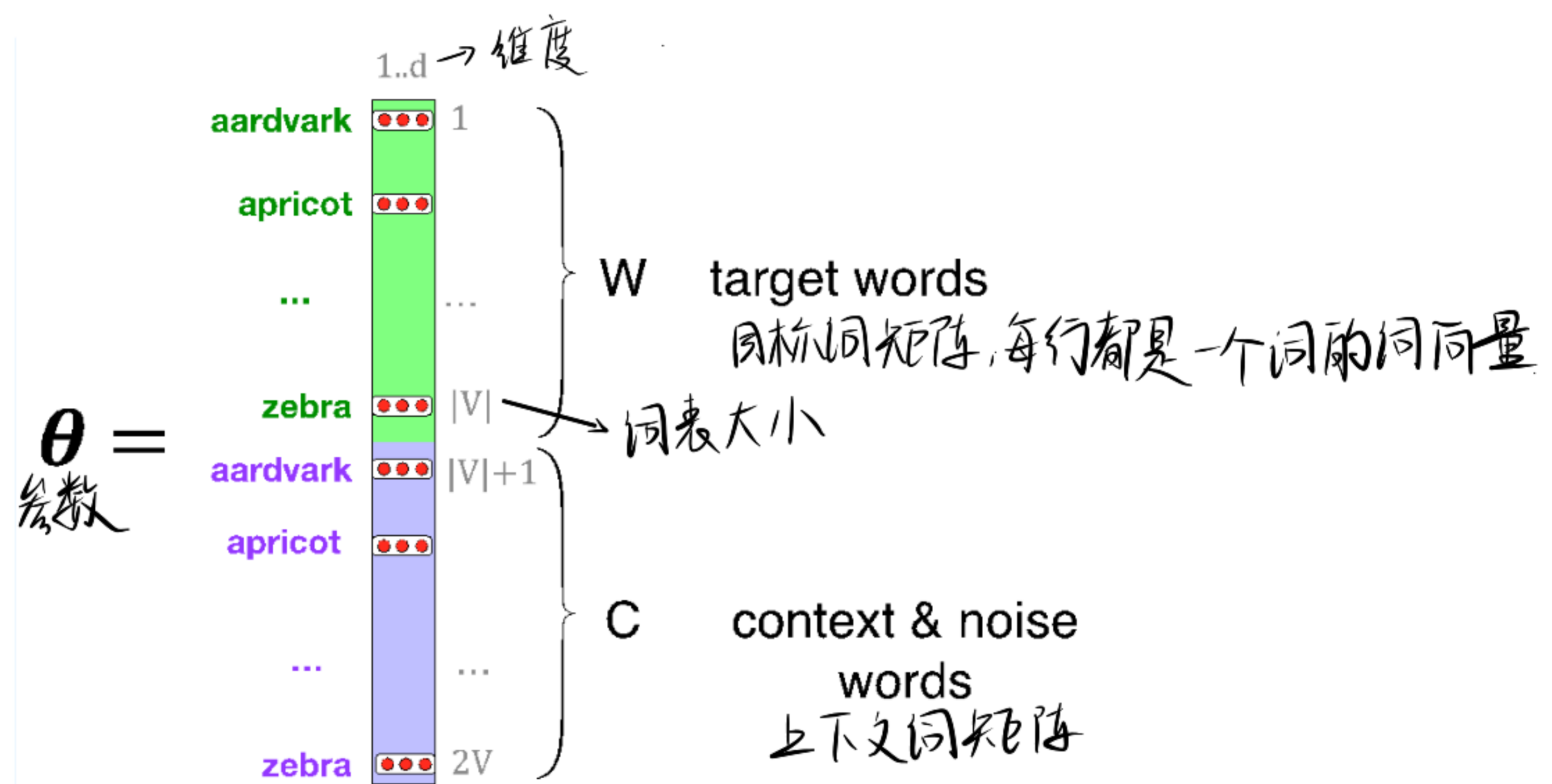
$$P(-|w, c) = 1 - P(+|w, c) = \frac{1}{1 + \exp(c \cdot w)}$$

假设所有上下文词语相互独立

$$P(+|w, c_{1:L}) = \prod_{i=1}^L \sigma(c_i \cdot w) \quad c_{1:L} \text{ 是窗口内 } L \text{ 个上下文词}$$

$$\text{对数化: } \log P(+|w, c_{1:L}) = \sum_{i=1}^L \log \sigma(c_i \cdot w)$$

接下来开始训练模型:



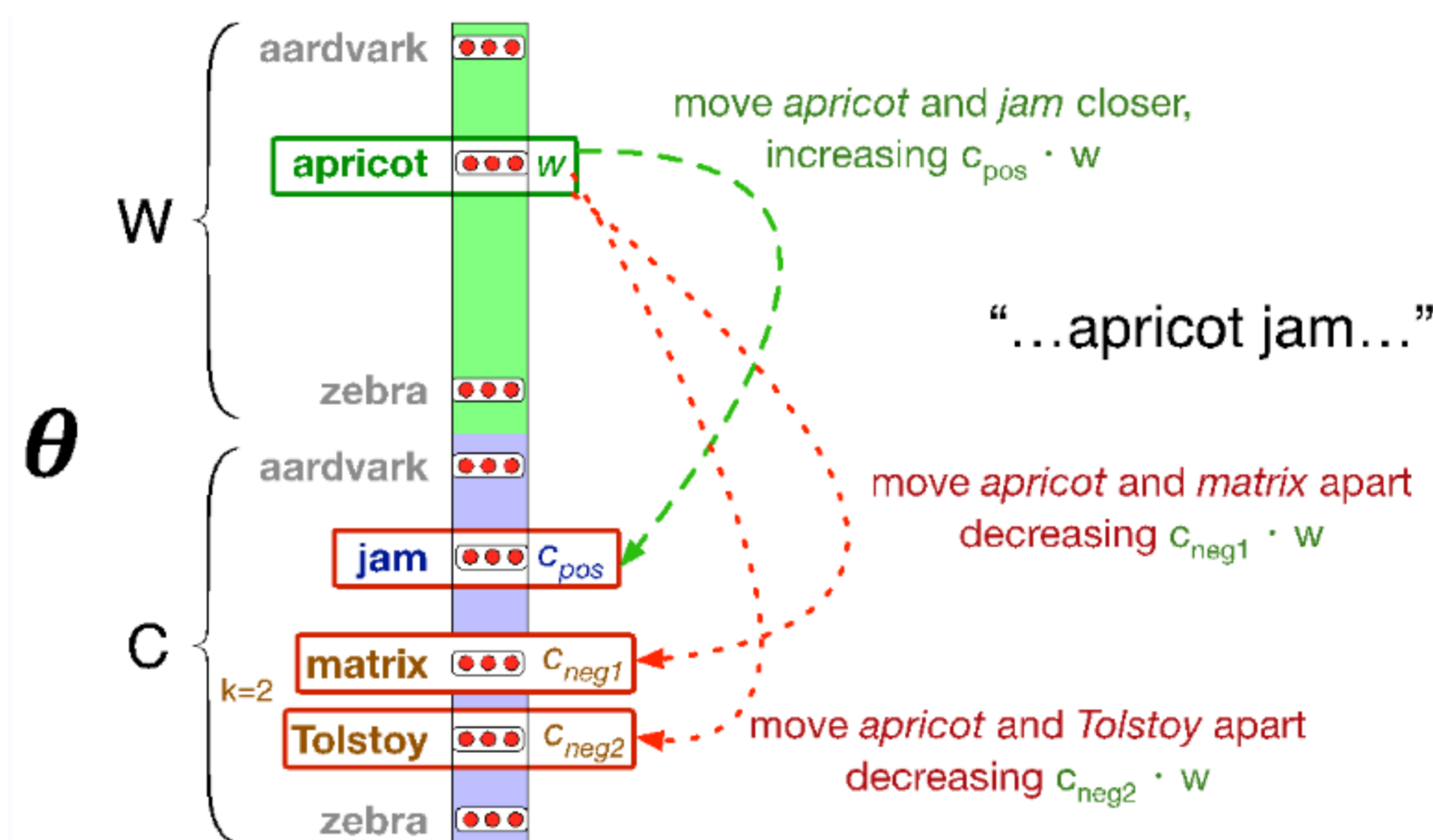
对于一个目标词, 我们选择一定数量的正样本, 对于每个正样本创建  $k$  个负样本

... lemon, a tablespoon of <b>apricot</b> jam a pinch				
...				
	c1	c2	t	c3 c4
<b>positive examples +</b>	<b>negative examples -</b> <span style="float: right;">k=2</span>			
t	c	t	c	t c
apricot	tablespoon	apricot	aardvark	apricot twelve
apricot	of	apricot	puddle	apricot hello
apricot	preserves	apricot	where	apricot dear
apricot	or	apricot	coaxial	apricot forever

给定正负训练样本和初始化的词嵌入向量 (即  $\theta$  中的  $W$  矩阵与  $C$  矩阵, 初始时所有值都是随机的), 不断调整

$W$  与  $C$  矩阵中的参数使得最大化  $P(W, C_{pos})$  相似度, 最小化  $P(W, C_{neg})$  相似度, 即最大化  $P(+|W, C_{pos})$  与  $P(-|W, C_{neg})$

$$\begin{aligned} \text{定义损失函数: } L_{CE} &= -\log \bar{P}(+|W, C_{pos}) \prod_{i=1}^k P(-|W, C_{neg_i}) \\ &= -\left[ \log P(+|W, C_{pos}) + \sum_{i=1}^k \log P(-|W, C_{neg_i}) \right] \\ &= -\left[ \log P(+|W, C_{pos}) + \sum_{i=1}^k \log (1 - P(+|W, C_{neg_i})) \right] \\ &= -\left[ \log \sigma(C_{pos} \cdot W) + \sum_{i=1}^k \log \sigma(-C_{neg_i} \cdot W) \right] \end{aligned}$$



通过随机梯度下降法调整词的权重(即词向量)最小化  $L_{CE}$

$$W^{t+1} = W^t - h \frac{d}{dW} L(f(x; W), y)$$

而对于这里的  $L_{CE}$ , 梯度为:

$$\frac{\partial L_{CE}}{\partial C_{pos}} = [\sigma(C_{pos} \cdot W) - 1] \cdot W$$

$$\frac{\partial L_{CE}}{\partial C_{neg}} = [\sigma(C_{neg} \cdot W)] \cdot W$$

$$\frac{\partial L_{CE}}{\partial W} = [\sigma(C_{pos} \cdot W) - 1] C_{pos} + \sum_{i=1}^k [\sigma(C_{neg_i} \cdot W)] C_{neg_i}$$

SGD更新方程:

$$C_{pos}^{t+1} = C_{pos}^t - \eta [\sigma(C_{pos}^t \cdot W^t) - 1] \cdot W^t$$

$$C_{neg}^{t+1} = C_{neg}^t - \eta [\sigma(C_{neg}^t \cdot W^t)] \cdot W^t$$

$$W^{t+1} = W^t - \eta \left\{ [\sigma(C_{pos}^t \cdot W^t) - 1] C_{pos}^t + \sum_{i=1}^k [\sigma(C_{neg_i}^t \cdot W^t)] C_{neg_i}^t \right\}$$

最后  $W$  和  $C$  都收敛, 得到每个目标词的权重, 作为该词的嵌入向量.

总结:

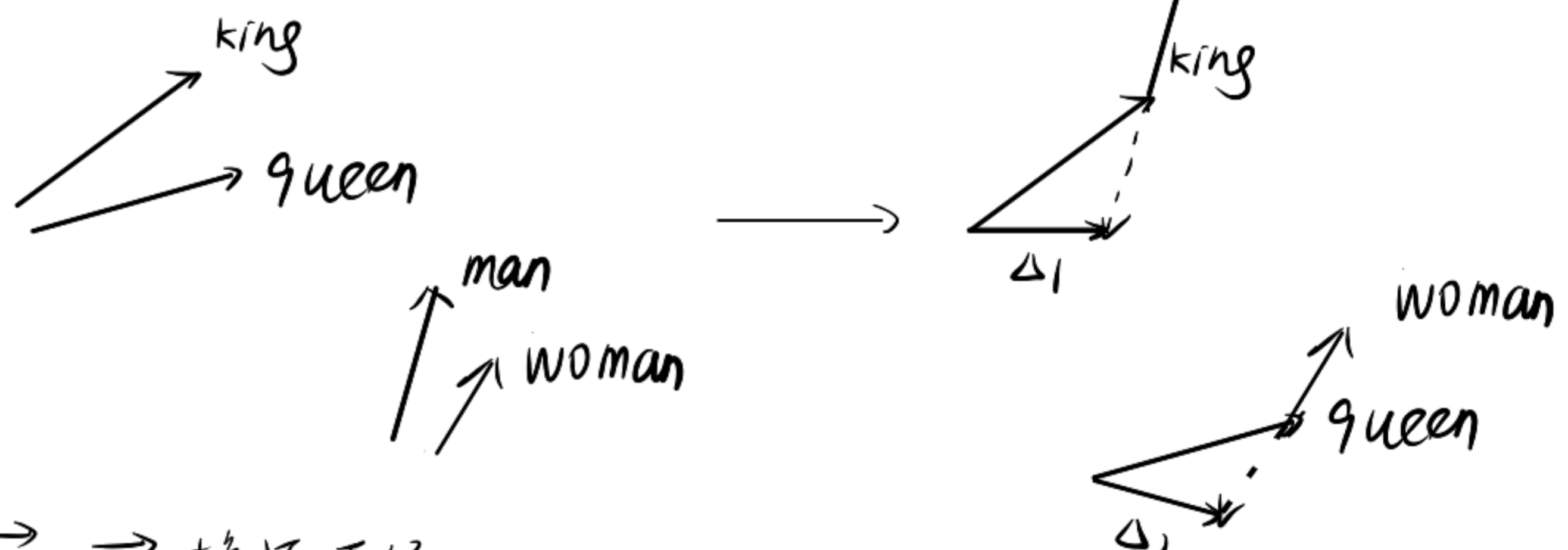
使用  $V$  个随机的  $d$ -维向量初始化嵌入向量  
基于嵌入向量相似度训练一个分类器  
从训练语料中选择上下文相关词对为正向样本  
选择训练语料中未出现的上下文相关词对为负向样本  
训练分类器同时调节所有词的嵌入向量直到损失函数达到最小  
分类器的权重向量为最终各个单词的嵌入向量

## Word2Vec 的推理能力

由于 Word2Vec 词向量是密集的、连续的, 它可以做一些运算, 而最有趣的是它总结规律、作出推理的能力

如:  $\text{king} - \text{man} + \text{woman} \approx \text{queen}$ .

数学语言:



$\vec{\Delta}_1, \vec{\Delta}_2$  接近平行, 即  $\text{man} - \text{king}$  与  $\text{woman} - \text{queen}$  偏移量几乎相同

这种能力使 Word2Vec 能学习复数、时态、比较级甚至风格、语气等规律.

而 TF-IDF 不具备这种能力, 因为它的向量维度太高且稀疏, 不成结构, 无法运算