

0 目录

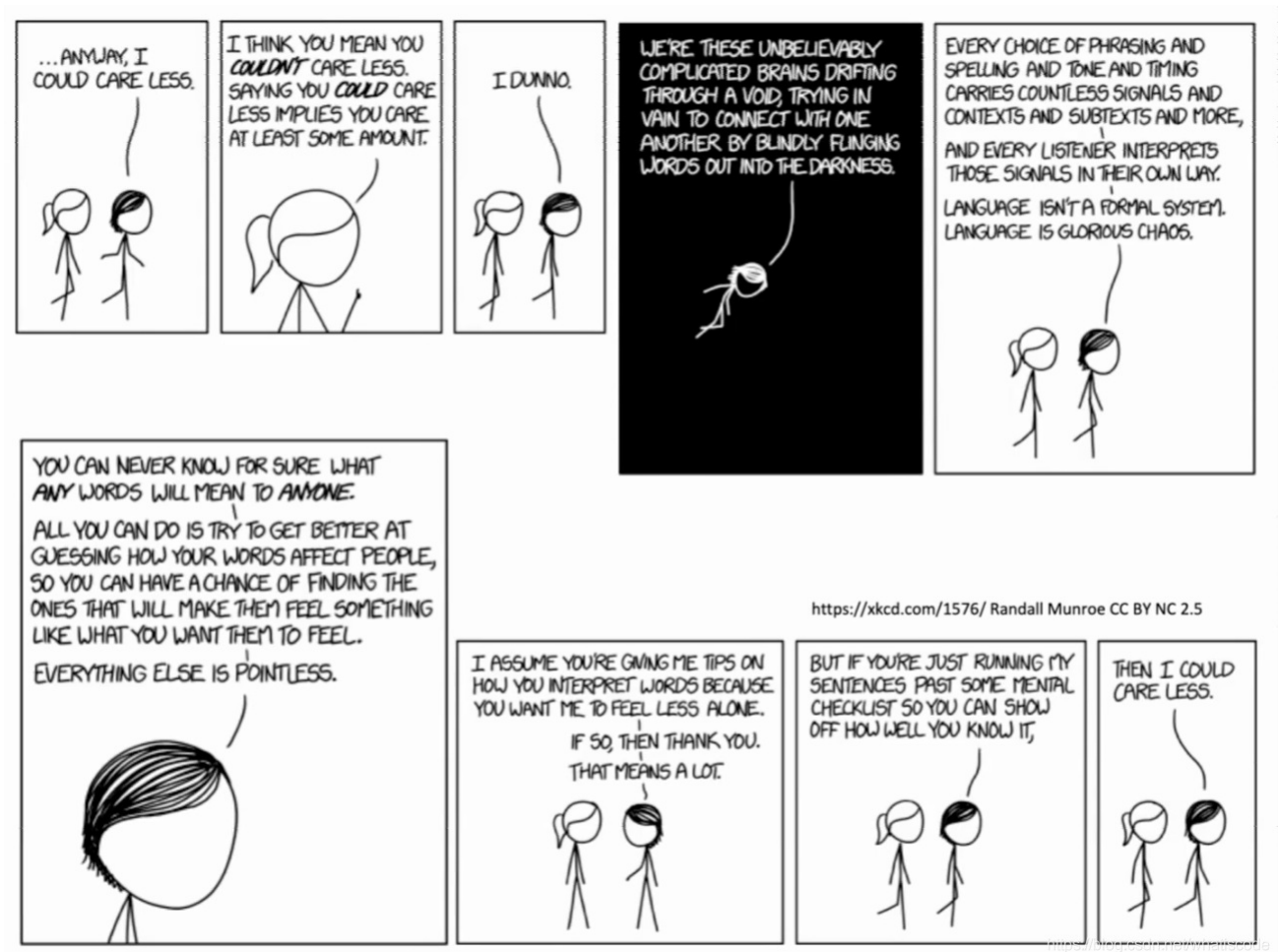
- 1. 教学目标
- 2. 人类语言和词义 (word meaning)
- 3. Word2vec介绍
- 4. Word2vec目标函数梯度
- 5. 优化方法

1 教学目标

- 1. 了解有效的现代深度学习方法
  - 首先学习基础知识，然后学习一些NLP领域中重要的方法：RNN，以及attention机制等等
- 2. 对人类语言的全局把握以及理解和形成人类语言的一些困难
- 3. 理解并掌握怎么去构造一个系统（PyTorch）来解决NLP领域中的一些主要的问题：
  - Word Meaning
  - Dependency Parsing
  - Machine Translation
  - Question Answering

2 人类语言和词义

2.1 一张XKCD漫画



<https://xkcd.com/1576/> Randall Munroe CC BY NC 2.5

<https://blog.csdn.net/whitewind/article/details/105444444>

- 语言带有不确定性
- 每一句话的措辞、拼写、语气以及说话的时间都带有无数的信号和上下文信息，并且每一个听众都以自己的方式去理解这些信号

2.2 人类语言

- 人类和大猩猩最大的区别就是人类拥有语言系统
- 人类大概在离开非洲时（约10万年前）拥有了语言
- 语言系统让人类能够更加有效的进行团队合作，这也是人类能够在自然界立于不败之地的原因
- 写作（writing）也是一项具有开创性的发明（5000年前）
- 写作让知识得以传承下去
- 语言系统能够起作用的原因是因为人类具有相似的知识储备。当我们用自然语言对一个场景进行描述时，听者在脑海中可以构建出相似的视觉场景。

2.3 我们怎么样表示一个词的意思（meaning）？

meaning的定义（韦伯斯特词典）：

- 用单词、词组表示概念
- 人们运用单词、符号表示自己的观点
- 通过写作作品、艺术来表达观点

signifier(symbol)⇔signified(idea or thing)

2.4 我们怎么在计算机上表示词义

2.4.1 WordNet

一个包含同义词（synonym）集合和上义词（hypernyms）的词典。

e.g. synonym sets containing “good”:

```
from nltk.corpus import wordnet as wn
poses = { 'n': 'noun', 'v': 'verb', 's': 'adj (s)', 'a': 'adj', 'r': 'adv' }
for synset in wn.synsets("good"):
    print("{}: {}".format(poses[synset.pos()],
        ", ".join([l.name() for l in synset.lemmas()])))
```

```
noun: good
noun: good, goodness
noun: good, goodness
noun: commodity, trade_good, good
adj: good
adj (sat): full, good
adj: good
adj (sat): estimable, good, honorable, respectable
adj (sat): beneficial, good
adj (sat): good
adj (sat): good, just, upright
...
adverb: well, good
adverb: thoroughly, soundly, good
```

e.g. hypernyms of “panda”:

```
from nltk.corpus import wordnet as wn
panda = wn.synset("panda.n.01")
hyper = lambda s: s.hypernyms()
list(panda.closure(hyper))
```

```
[Synset('procyonid.n.01'),
Synset('carnivore.n.01'),
Synset('placental.n.01'),
Synset('mammal.n.01'),
Synset('vertebrate.n.01'),
Synset('chordate.n.01'),
Synset('animal.n.01'),
Synset('organism.n.01'),
Synset('living_thing.n.01'),
Synset('whole.n.02'),
Synset('object.n.01'),
Synset('physical_entity.n.01'),
Synset('entity.n.01')]
```

<https://blog.csdn.net/whatiscode>

Wordnet存在的问题:

- 作为一个较好的资料，但忽略了细微的一些差别：例如词典中‘proficient’与‘good’认为是同义词，但是这只在某些文本上下文中成立。
- 忽略了一些单词的新的含义（无法随时更新）
- 比较主观（缺少客观性）
- 需要人类劳动力来不断地创建和更新
- 不能计算单词之间的相似度

#### 2.4.2 Representing words as discrete symbols

传统NLP中，我们将单词看作是离散的表示，通过one-hot vector来表示：

Means one 1, the rest 0s

Words can be represented by one-hot vectors:

```
motel = [0 0 0 0 0 0 0 0 0 0 1 0 0]
hotel = [0 0 0 0 0 0 0 1 0 0 0 0 0]
```

One-hot表示存在的问题:

- 如果词汇量太多，向量的维度会变得很大。
- 由于向量是正交的，因此没法表示他们之间啊的相似度。

#### 2.4.3 Representing words by their context

Distributional semantics: 一个单词的含义通常由在它附近经常出现的单词给出的。

对于文本中的一个单词w，它的上下文就是出现在它附近的一组单词（在一个划定好size的窗口下）

通过许多包含w的文本中的上下文来构建w的含义表示：

```
...government debt problems turning into banking crises as happened in 2009...
...saying that Europe needs unified banking regulation to replace the hodgepodge...
...India has just given its banking system a shot in the arm...
```

These context words will represent **banking**

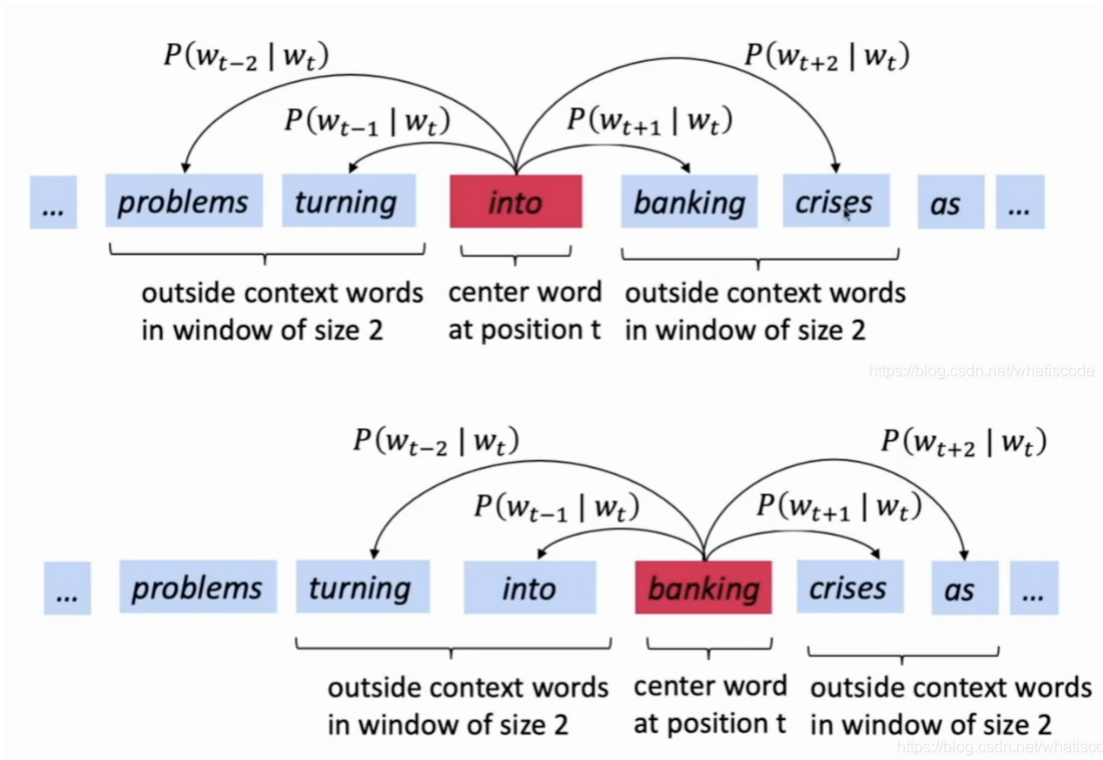
3 Word2vec介绍

Word2vec是一种学习词向量的框架。

3.1主要思想

- 包含大量的文本语料
- 固定词表中的每一个单词由一个词向量表示
- 文本中的每个单词位置  $t$ ，有一个中心词  $c$ ，和它的上下文  $o$ （除了  $c$  的外部单词）。
- 通过  $c$  和  $o$  的词向量相似性来计算  $P(o|c)$
- 不断的调整词向量，最大化概率

3.2 计算  $P(w_{t+j} | w_t) P(w_{-t+j}|w_t) P(w_{t+j} | w_t)$  的示例



3.3 Word2vec的目标函数

对于每个位置  $t=1,...,T$ ，固定窗口大小  $m$ ，给定中心词  $w_j w_{-j} w_j$ :

$likelihood = L(\theta) = \prod_{t=1}^T \prod_{-m \leq j \leq m, j \neq 0} p(w_{t+j} | w_t; \theta)$   
 $likelihood = L(\theta) = \prod_{t=1}^T \prod_{-m \leq j \leq m, j \neq 0} p(w_{t+j} | w_t; \theta)$

- $\theta$  是需要优化的参数

$J(\theta) = -1/T \log L(\theta) = -1/T \sum_{t=1}^T \sum_{-m \leq j \leq m, j \neq 0} \log p(w_{t+j} | w_t; \theta)$   
 $J(\theta) = -1/T \log L(\theta) = -1/T \sum_{t=1}^T \sum_{-m \leq j \leq m, j \neq 0} \log p(w_{t+j} | w_t; \theta)$

- $J(\theta)$  为损失函数（这里是平均负对数似然）；
- 负号将最大化损失函数转化为最小化损失函数；
- $\log$  函数方便将乘法转化为求和（优化处理）

3.4 如何计算  $P(w_{t+j} | w_t; \theta) P(w_{-t+j}|w_t; \theta) P(w_{t+j} | w_t; \theta)$

对于每个单词  $w$  我们使用两个向量  $v_w v_w$  和  $u_w u_w$

- $v_w v_w$  表示当  $w$  是中心词时
- $u_w u_w$  表示当  $w$  是上下文单词时
- 对于中心词  $c$  和上下文单词  $o$ ，有： $P(o | c) = \frac{\exp(u_o^T v_c)}{\sum_{w \in V} \exp(u_w^T v_c)}$   $P(o | c) = \frac{\exp(u_o^T v_c)}{\sum_{w \in V} \exp(u_w^T v_c)}$

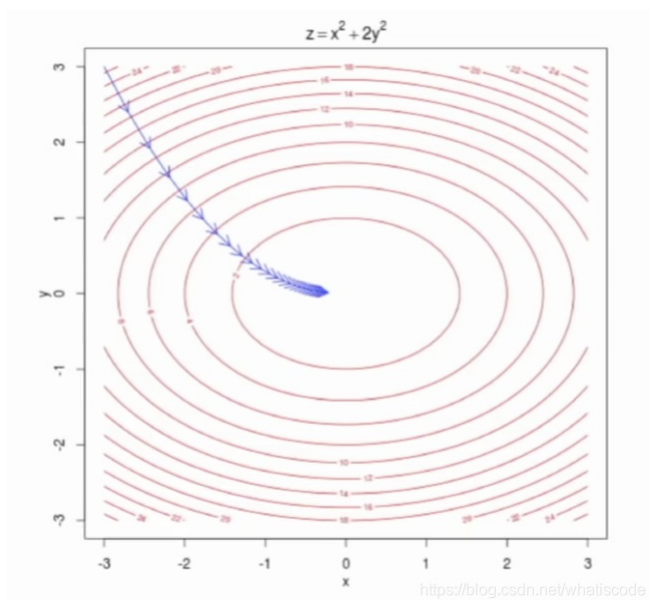
在上式中：

- 分子取幂函数使得始终可以为正
- 向量  $u_o u_o$  和向量  $v_c v_c$  点乘，点乘结果越大，向量之间越相似
- $u^T v = u \cdot v = \sum_i u_i v_i$   $u^T v = u \cdot v = \sum_i u_i v_i$
- 对整个词表标准化，给出概率分布
- softmax函数进行归一化（深度学习常用）： $R^n \rightarrow R^n$   $\text{softmax}(x) = \frac{\exp(x_i)}{\sum_j \exp(x_j)}$   $\text{softmax}(x) = \frac{\exp(x_i)}{\sum_j \exp(x_j)}$

4 Word2vec目标函数梯度

4.1 训练模型

通过调整参数的方式来最小化损失函数



#### 4.1.1 训练模型的方法：计算所有的向量梯度

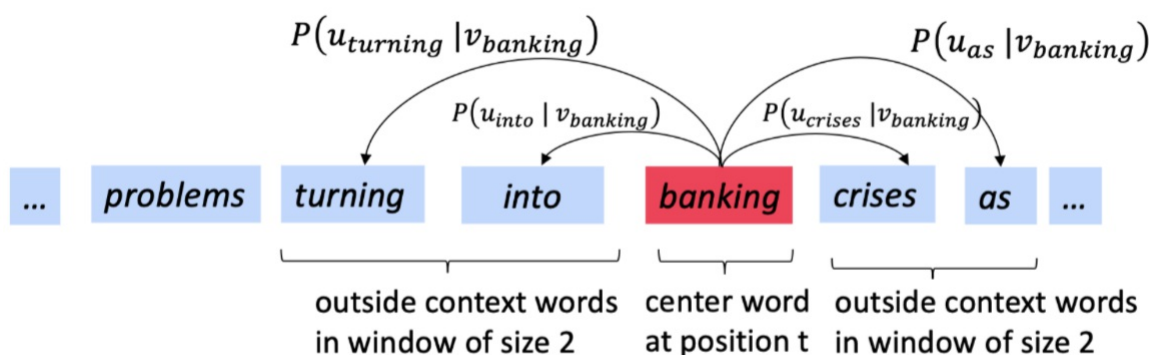
由于整个模型里只有一个参数  $\theta$ , 所以我们只要优化这一个参数就行。如一个  $d$  维，词典大小为  $V$  的模型所包含的参数（每个单词有两个向量）：

$$\theta = \begin{bmatrix} v_{aardvark} \\ v_a \\ \vdots \\ v_{zebra} \\ u_{aardvark} \\ u_a \\ \vdots \\ u_{zebra} \end{bmatrix} \in \mathbb{R}^{2dV}$$

我们可以通过梯度下降的方式优化参数，梯度下降会用到链式法则。

- 迭代计算每个中心词向量和上下文词向量随着滑动窗口移动的梯度
- 依次迭代更新窗口中所有的参数

示例：



## 5 优化方法

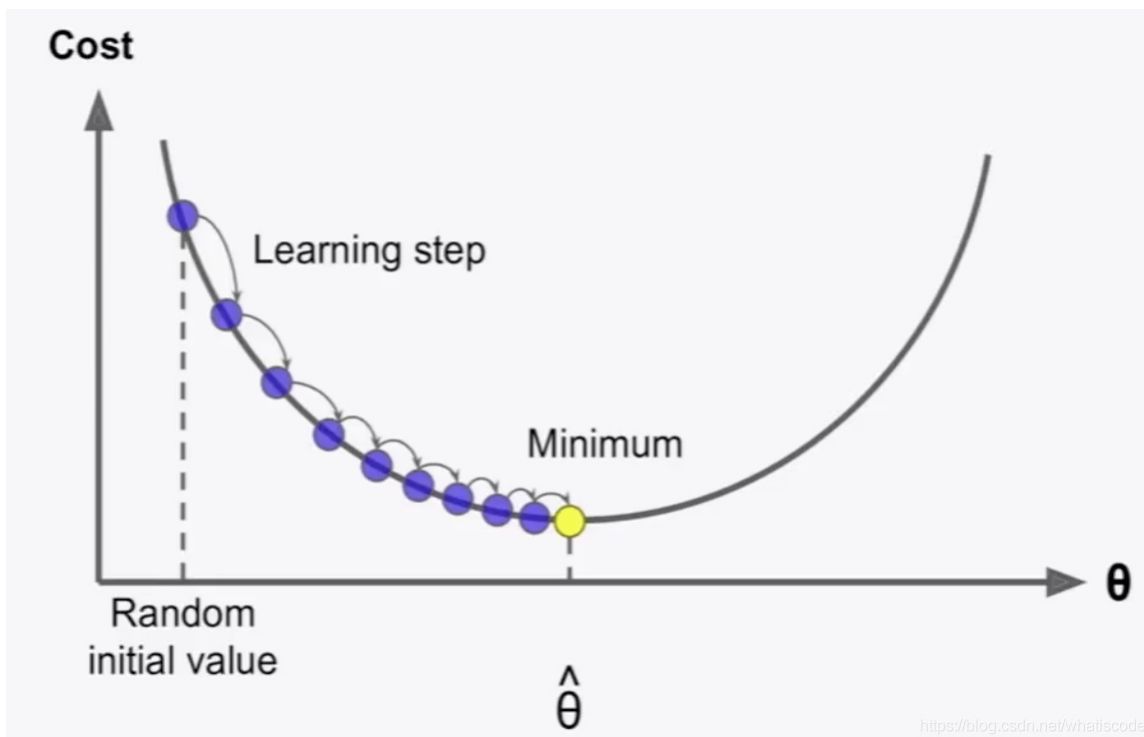
### 5.1 梯度下降法

最小化损失函数  $J(\theta)$

对于当前  $\theta$ , 计算  $J(\theta)$  的梯度

然后小步重复朝着负梯度方向更新方程里的参数  $\alpha = (\text{step size}) \text{ or } (\text{learning rate})$   
 $\theta_{\text{new}} = \theta_{\text{old}} - \alpha \nabla J(\theta)$

更新唯一的参数  $\theta_{\text{new}} = \theta_{\text{old}} - \alpha \nabla J(\theta)$



```
while True:
    theta_grad = evaluate_gradient(J,corpus,theta)
    theta = theta - alpha * theta_grad
```

## 5.2 随机梯度下降SGD

由于  $J(\theta)$  是在语料文本中所有窗口的方程

当语料很大的时候，计算梯度会消耗巨大

解决办法: SGD

不断sample窗口，不断更新

```
while True:
    window = sample_window(corpus)
    theta_grad = evaluate_gradient(J>window,theta)
    theta = theta - alpha * theta_grad
```