

开始：文本表示

Word Representation

One-hot
representation

词典: [我们, 去, 爬山, 今天, 你们, 昨天, 跑步]

①

③

⑥

⑦

每个单词的表示:

- 我们: $(1, 0, 0, 0, 0, 0, 0) \rightarrow 7\text{维} = |\text{词典}|$
- 爬山: $(0, 0, 1, 0, 0, 0, 0) \rightarrow 7\text{维} = |\text{词典}|$
- 跑步: $(0, 0, 0, 0, 0, 0, 1) \rightarrow 7\text{维} = |\text{词典}|$
- 昨天: $(0, 0, 0, 0, 0, 1, 0) \rightarrow 7\text{维} = \cdot\cdot\cdot$

Sentence Representation (boolean)

Boolean
Representation

词典: [我们, 又, 去, 爬山, 今天, 你们, 昨天, 跑步]

8个单词

每个句子的表示

- { 我们 今天 去 爬山: $(1, 0, 1, 1, 1, 0, 0, 0) \rightarrow 8\{\} = 12$ |
你们 昨天 跑步: $(0, 0, 0, 0, 0, 1, 1, 1) \rightarrow 8\{\} = 12$ |
你们 又 去 爬山 又 去 跑步: $(0, 1, 1, 1, 0, 1, 0, 1) \rightarrow 8\{\} = 12$ |

Sentence Representation (count)

Count-based
representation



词典: [我们, 又, 去, 爬山, 今天, 你们, 昨天, 跑步]

每个句子的表示

我们 今天 去 爬山: $(1, 0, 1, 1, 1, 0, 0, 0) \Rightarrow 8\{3 = 1\} \text{ 是}$

你们 昨天 跑步: $(0, 0, 0, 0, 0, 1, 1, 1) \Rightarrow 8\{3 = 1\} \text{ 是}$

你们 又 去 爬山 又 去 跑步: $(0, 2, 2, 1, 0, 1, 0, 1) \Rightarrow 8\{3 = 1\} \text{ 是}$

Sentence Similarity

Distance

计算距离 (欧式距离) : $d = |s1 - s2|$

S1: “我们今天去爬山” = (1, 0, 1, 1, 0, 0, 0, 0)

S2: “你们昨天跑步” = (0, 0, 0, 0, 0, 1, 1, 1)

S3: “你们又去爬山又去跑步” = (0, 2, 2, 1, 0, 1, 0, 1)

$$d(s_1, s_2) = \sqrt{1^2 + 1^2 + 1^2 + 1^2 + 1^2 + 1^2} = \sqrt{6}$$

$$d(s_1, s_3) = \sqrt{1^2 + 2^2 + 1^2 + 1^2 + 1^2} = \sqrt{8}$$

$$d(s_2, s_3) = \sqrt{2^2 + 2^2 + 1^2 + 1^2} = \sqrt{10}$$

$$\text{sim}(s_1, s_2) > \text{sim}(s_2, s_3); \quad \text{sim}(s_1, s_3) > \text{sim}(s_2, s_3)$$

$$s_1 = (x_1, x_2, x_3)$$

$$s_2 = (y_1, y_2, y_3)$$

$$d = \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2 + (x_3 - y_3)^2}$$

Sentence Similarity

计算相似度（余弦相似度）： $d = \frac{\langle s_1, s_2 \rangle}{\|s_1\| * \|s_2\|}$

内积

Normaliz...

S1: “我们今天去爬山” = (1,0,1,1,0,0,0,0)

S2: “你们昨天跑步” = (0,0,0,0,1,1,1)

S3: “你们又去爬山又去跑步” = (0,2,2,1,0,1,0,1)

$$d(s_1, s_2) = 0/A = 0$$

$$d(s_1, s_3) = \frac{(2+1)}{\sqrt{3} \cdot \sqrt{11}} = \frac{3}{\sqrt{33}}$$

$$d(s_2, s_3) = \frac{2}{\sqrt{3} \cdot \sqrt{11}} = \frac{2}{\sqrt{33}}$$

$$\begin{aligned}s_1 &= (x_1, x_2, x_3) \\s_2 &= (y_1, y_2, y_3)\end{aligned}$$

$$d = \frac{x_1y_1 + x_2y_2 + x_3y_3}{\sqrt{x_1^2 + x_2^2 + x_3^2} \sqrt{y_1^2 + y_2^2 + y_3^2}}$$

$$\text{sim}(s_1, s_3) > \text{sim}(s_2, s_3) > \text{sim}(s_1, s_2)$$

Cosine Similarity

Sentence Similarity

句子1: He is going from Beijing to Shanghai

句子2: He denied my request, but he actually lied.

句子3: Mike lost the phone, and phone was in the car

=> i.e. 哪

句子1: (0, 0, 1, 0, 0, 0, 1, 1, 1, 0, 1, 0, 0, 0, 0, 0, 1, 0, 1, 0)

句子2: (1, 0, 0, 1, 0, 1, 0, 0, 2, 0, 0, 1, 0, 0, 1, 0, 1, 0, 0, 0)

句子3: (0, 1, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 1, 1, 0, 2, 0, 0, 2, 0, 1)

count-based
representation

denied > he

Sentence Similarity

句子1: He is going from Beijing to Shanghai

句子2: He denied my request, but he actually lied.

句子3: Mike lost the phone, and phone was in the car

句子1: (0, 0, 1, 0, 0, 0, 1, 1, 1, 0, 1, 0, 0, 0, 0, 0, 1, 0, 1, 0)

句子2: (1, 0, 0, 1, 0, 1, 0, 0, 2, 0, 0, 1, 0, 0, 1, 0, 1, 0, 0, 0, 0)

句子3: (0, 1, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 1, 1, 0, 2, 0, 0, 2, 0, 1)

denied

he

① 并不是出现的越多就越重要
② 并不是出现的越少就越不重要!

Sentence Similarity

句子1: He is going from Beijing to Shanghai

句子2: He denied my request, but he actually lied.

句子3: Mike lost the phone, and phone was in the car

句子1: (0, 0, 1, 0, 0, 0, 1, 1, 1, 0, 1, 0, 0, 0, 0, 0, 1, 0, 1, 0)

句子2: (1, 0, 0, 1, 0, 1, 0, 0, 2, 0, 0, 1, 0, 0, 1, 0, 1, 0, 0, 0, 0)

句子3: (0, 1, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 1, 1, 0, 2, 0, 0, 2, 0, 1)

denied

he

{ 并不是出现的越多就越重要!
 { 并不是出现的越少就越不重要!

Tf-idf Representation

$$tfidf(w) = tf(d, w) * idf(w) \rightarrow \text{考虑单词的重要性}$$

文档 d 中 w 的词频

$$\log \frac{N}{N(w)}$$

i.e., $N=100$
 $N(w)=5$
 $\frac{N}{N(w)} = 20$

N: 语料库中的文档总数

$N(w)$: 词语 w 出现在多少个文档?

$tfidf(w) = tf(d, w) * idf(w)$ ① 词典: [今天, 上, NLP, 课程, 的, 有, 意思, 数据, 也]
|词典| = 9

句1 今天 上 NLP 课程

$\Rightarrow \underline{\text{句1}}$

$$= (1 \cdot \log \frac{3}{2}, 1 \cdot \log \frac{3}{1}, 1 \cdot \log \frac{3}{1}, 1 \cdot \log \frac{3}{3}, 0, 0, 0, 0, 0)$$

$$= (\underbrace{\log \frac{3}{2}, \log 3, \log 3, \log 1, 0, 0, 0, 0, 0}_{tf-idf \text{ 向量}}) \Rightarrow 9 组 = |\text{词典}|$$

$\Rightarrow \underline{\text{句2}}$

$$= (1 \cdot \log \frac{3}{2}, 0, 0, 1 \cdot \log \frac{3}{3}, 1 \cdot \log \frac{3}{1}, 1 \cdot \log \frac{3}{2}, 1 \cdot \log \frac{3}{2}, 0, 0)$$

$$= (\underbrace{\log \frac{3}{2}, 0, 0, \log 1, \log 3, \log \frac{3}{2}, \log \frac{3}{2}, 0, 0}_{tf-idf \text{ 向量}}) \Rightarrow 9 组 = |\text{词典}|$$

$\Rightarrow \underline{\text{句3}}$

$$= (0, 0, 0, 1 \cdot \log \frac{3}{3}, 0, 1 \cdot \log \frac{3}{2}, 1 \cdot \log \frac{3}{2}, 1 \cdot \log \frac{3}{1}, 1 \cdot \log \frac{3}{1})$$

$$= (0, 0, 0, \log 1, 0, \log \frac{3}{2}, \log \frac{3}{2}, \log 3, \log 3) \Rightarrow 9 组 = |\text{词典}|$$

Measure Similarity Between Words

下面哪些单词之间语义相似度更高？

我们，爬山，运动，昨天

$$\text{我们} = (0, 1, 0, 0, \dots)$$

$$\text{爬山} = (0, 0, 1, 0, \dots)$$

$$\text{sim}(\text{我们}, \text{爬山}) < \text{sim}(\text{运动}, \text{爬山})$$

- o One-hot representation
 - boolean-based representation
 - count-based "
 - tf-idf-based representation

Measure Similarity Between Words

下面哪些单词之间语义相似度更高？

我们，爬山，运动，昨天

我们：(0, 1, 0, 0, 0, 0)
爬山：(0, 0, 1, 0, 0, 0)
运动：(1, 0, 0, 0, 0, 0)
昨天：(0, 0, 0, 1, 0, 0)

① Euclidean Distance

$$d(\text{我们}, \text{爬山}) = \sqrt{2}$$

$$d(\text{我们}, \text{运动}) = \sqrt{2}$$

$$d(\text{运动}, \text{爬山}) = \sqrt{2}$$

$$d(\text{运动}, \text{昨天}) = \sqrt{2}$$

X

② Cosine Similarity
 $\text{sim}(\text{我们}, \text{爬山}) = \frac{0}{\sqrt{2}} = 0$
 $\text{sim}(\text{我们}, \text{运动}) = 0$
 $\text{sim}(\text{运动}, \text{爬山}) = 0$
 $\text{sim}(\text{运动}, \text{昨天}) = 0$

X

One-hot representation
有可能(用)来表示语义相似度？

Measure Similarity Between Words

利用 One-hot 表示法表达单词之间相似度?

每个单词的表示:

我们: [1, 0, 0, 0, 0, 0, 0]

爬山: [0, 0, 1, 0, 0, 0, 0]

运动: [0, 0, 0, 0, 0, 0, 1]

昨天: [0, 0, 0, 0, 0, 1, 0]



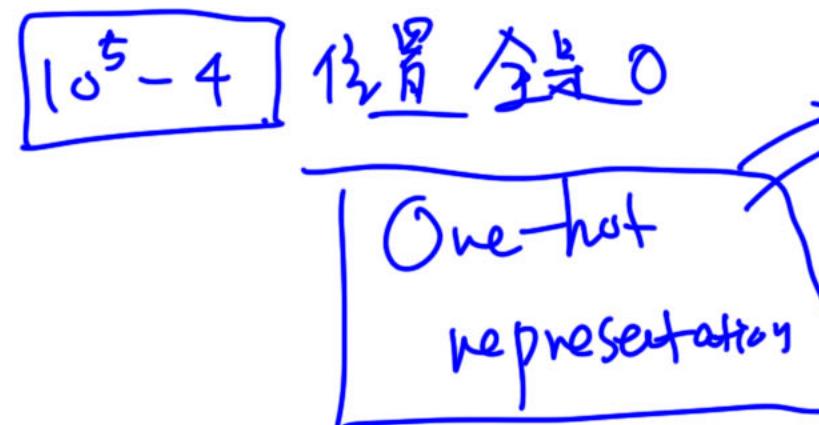
N G ↗
↓

Another Issue: Sparsity

1 我们今天打算去爬山 = () \Rightarrow 向量的大小 = | i_2 | 哪

2 你们昨天做什么了 \Rightarrow

3 明天打算去上课 \Rightarrow $(\underbrace{0.0 \dots 1.0 \dots 0.1 \dots 1.1}_{10^5 \sim 10^6})$ i.e., 中文 i_2 哪
 $10^5 \sim 10^6$



① 不能表示语义
的相似性

② Sparsity

From One-hot Representation to Distributed Representation

分布式表示法？

One-Hot Representation

我们: [1, 0, 0, 0, 0, 0]

爬山: [0, 0, 1, 0, 0, 0]

运动: [0, 0, 0, 0, 0, 1]

昨天: [0, 0, 0, 0, 1, 0]



Distributed Representation •

我们: [0.1, 0.2, 0.4, 0.2]

爬山: [0.2, 0.3, 0.7, 0.1]

运动: [0.2, 0.3, 0.6, 0.2]

昨天: [0.5, 0.9, 0.1, 0.3]

长度 = 1词典
 $10^5 \sim 10^6$



长度 = 100, 200, 300
(Sparsity)

Measure Similarity Between Words

Distributed Representation

我们: [0.1, 0.2, 0.4, 0.2]

爬山: [0.2, 0.3, 0.7, 0.1]

运动: [0.2, 0.3, 0.6, 0.2]

昨天: [0.5, 0.9, 0.1, 0.3]

四维向量

① Euclidean Distance

$$d(\text{我们}, \text{爬山}) = \sqrt{0.1^2 + 0.1^2 + 0.3^2 + 0.1^2} \\ = \sqrt{0.01 + 0.01 + 0.09 + 0.01} = \sqrt{0.12}$$

$$d(\text{运动}, \text{爬山}) = \sqrt{0.1^2 + 0.1^2} = \sqrt{0.02}$$

$$(d(\text{运动}, \text{爬山}) < d(\text{我们}, \text{爬山}))$$

$$\Rightarrow \text{sim}(\text{运动}, \text{爬山}) > \text{sim}(\text{我们}, \text{爬山})$$

$$② d(\text{爬山}, \text{运动}) < d(\text{昨天}, \text{运动})$$

$$\Rightarrow \text{sim}(\text{爬山}, \text{运动}) > \text{sim}(\text{昨天}, \text{运动})$$

分布式表示方法(单词)

词向量(word vectors)

Comparing the Capacities

Q: 100 维的 One-Hot 表示法最多可以表达多少个不同的单词?

Q: 100 维的 分布式 表示法最多可以表达多少个不同的单词?

Comparing the Capacities

容量空间)

Q: 100 维的 One-Hot 表示法最多可以表达多少个不同的单词?

向量大小 = 100

我们 = (1, 0, 0...0)
 99个

100 个单词 !

运动 = (0, 0, 0...1, 0...0)

Q: 100 维的 分布式 表示法最多可以表达多少个不同的单词?

我们: (0.1, 0.2, 0.15...0.1)
 100 维

+∞

Binary

我们: (1, 0, 1, 0...0, 1...)
 100 位
 ↓ ↓
 0/1 0/1

2^{100} 不同的单词 !

Questions

Q: 怎么学习每一个单词的分布式表示（词向量）？

我们： $(0.2, 0.3, 0.1, 0.2)$ ？

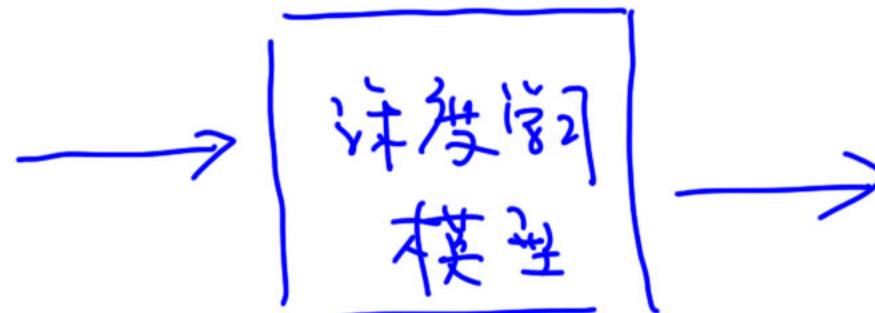
Learn Word Embeddings

输入 (Input) : String

我们今天去爬山

你么昨天运动

你们去爬山



1B: string contains 10 tokens

10B

100B

Distributed Representation

我们: [0.1, 0.2, 0.4, 0.2]

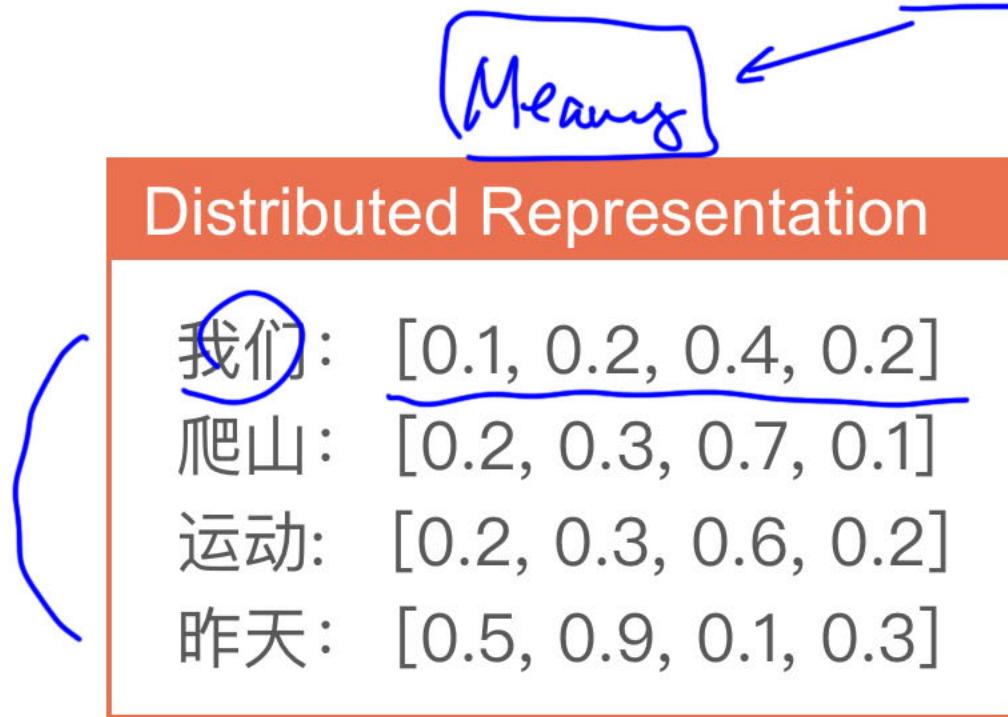
爬山: [0.2, 0.3, 0.7, 0.1]

运动: [0.2, 0.3, 0.6, 0.2]

昨天: [0.5, 0.9, 0.1, 0.3]

dim/D: 100/200/300, 50 *

Essence of Word Embedding

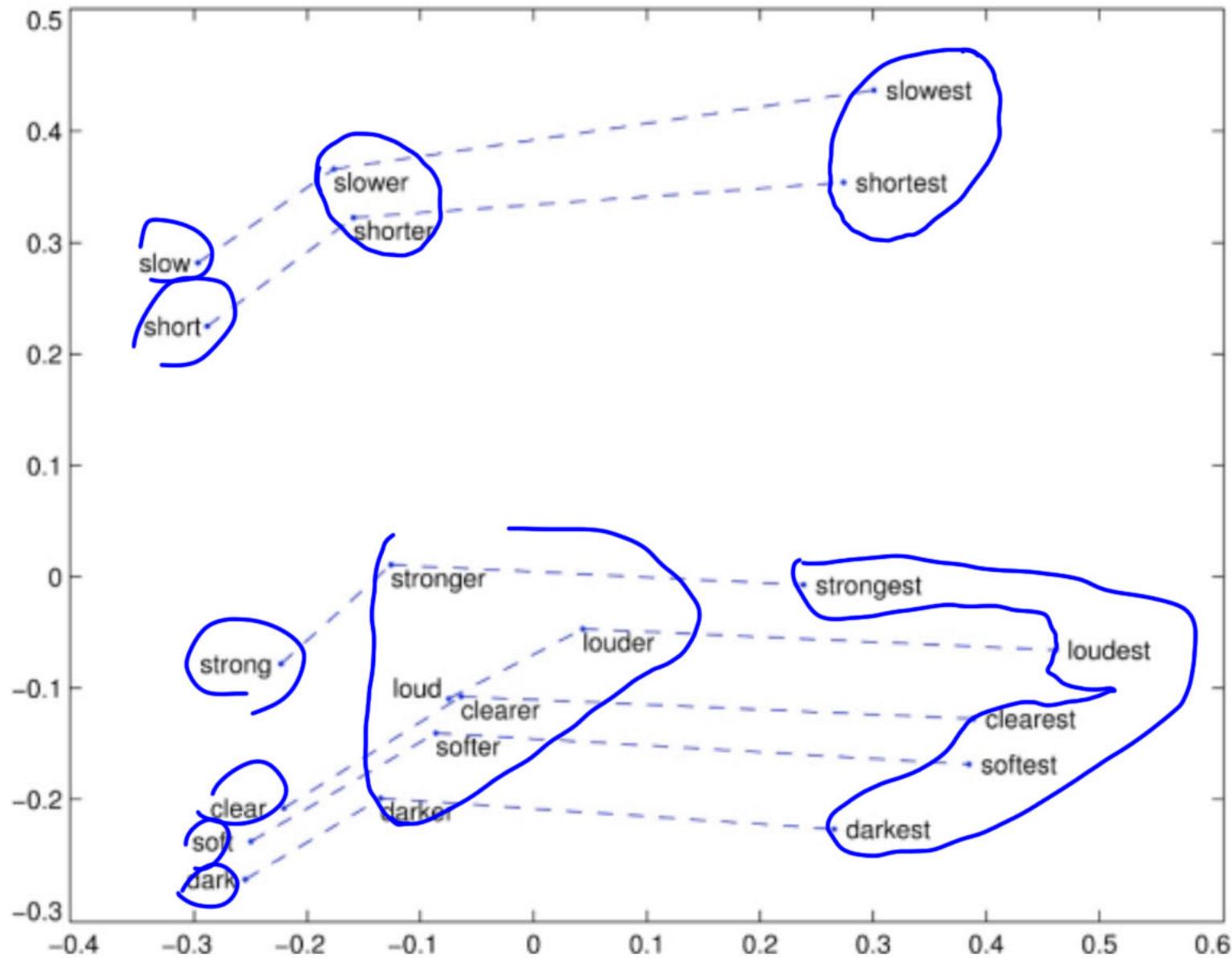


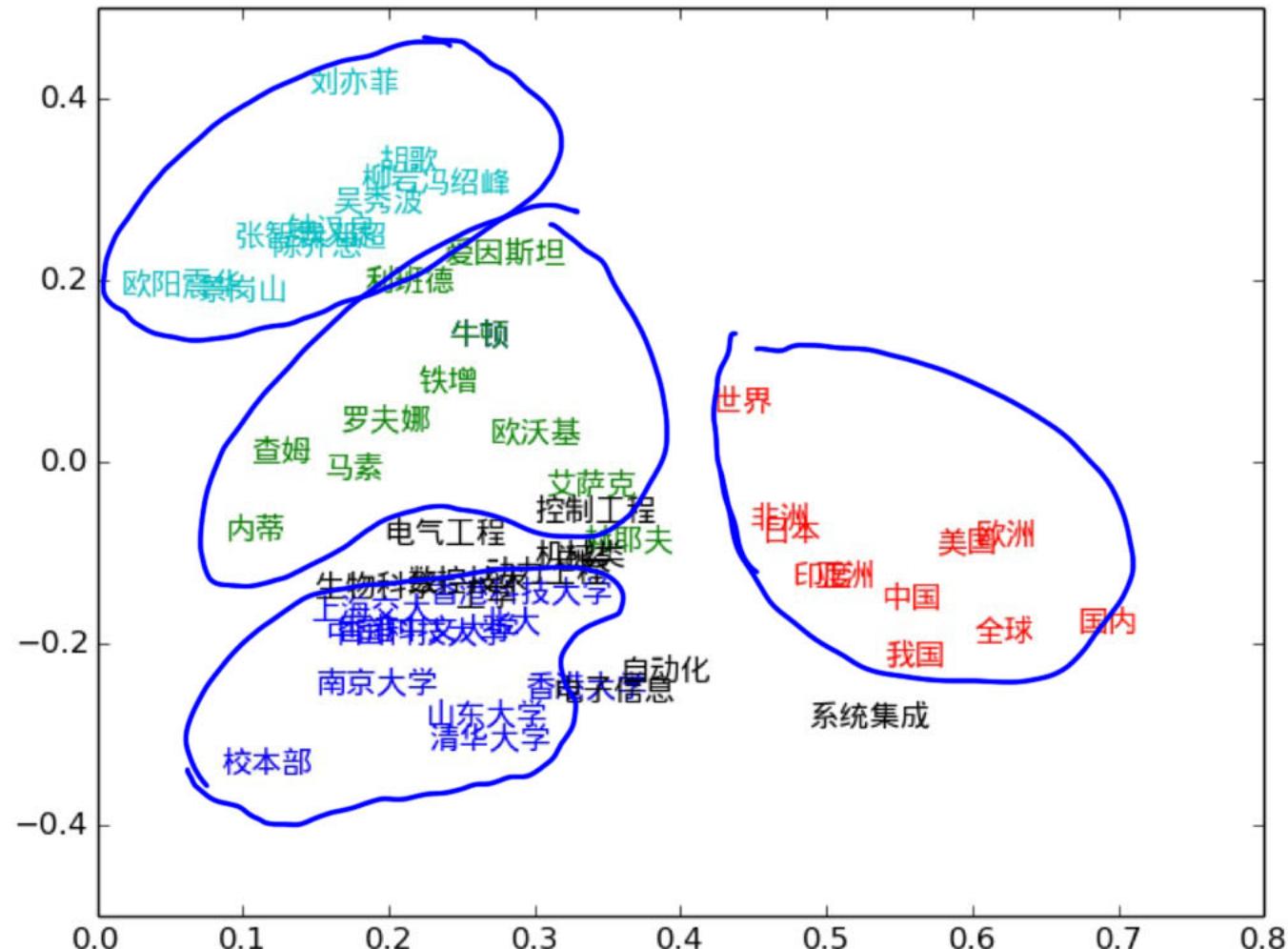
woman - man
≈ girl + boy



词向量代表单词的意义
(Meaning)

Word2Vec 萤幕上
矩阵式词向量
(Meaning)





From Word Embedding to Sentence Embedding

我们 : $(0.1, 0.2, 0.1, 0.3)$

去 : $(0.3, 0.2, 0.15, 0.2)$

+ 运动 : $(0.2, 0.15, 0.4, 0.7)$

0.6 0.55 0.65 1.2

$(0.2 0.18 0.22 0.4)$

A. “我们 去 运动”

① Averaging 方法

句子向量

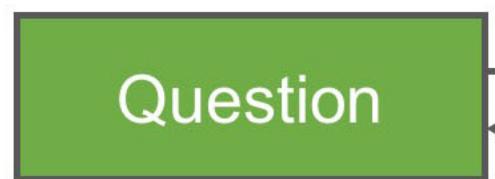
$= (0.2, 0.18, 0.22, 0.4) \checkmark$

B: “我们去爬山” $= (0.2, 0.2, 0.25, 0.4) \checkmark$

② LSTM/RNN

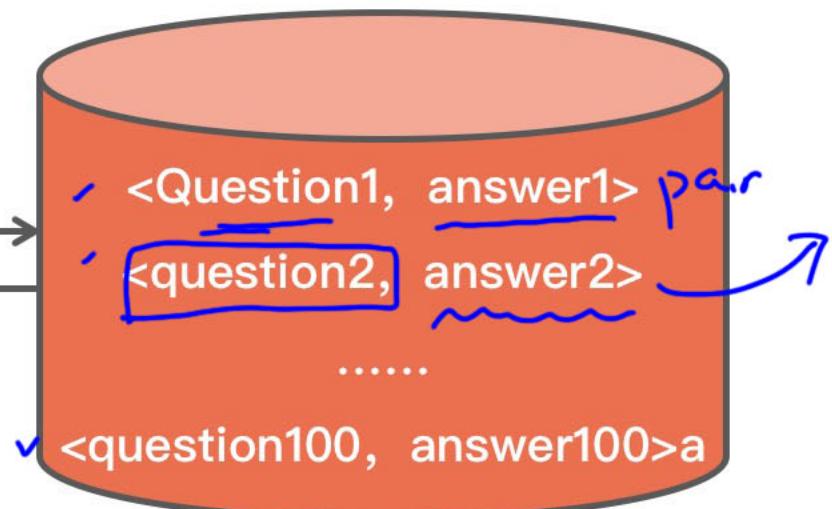
Recap: Retrieval-based QA System

How do you like NLPCamp?



相似度匹配

返回相似度最高的

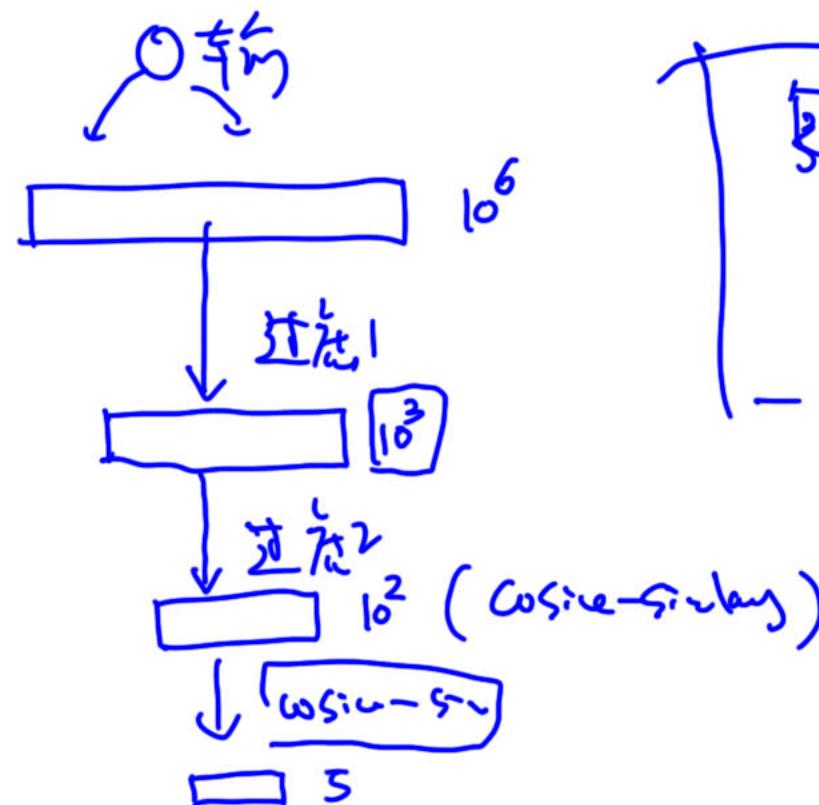


$O(N)$ · 每次相似度计算复杂度
 $\sin(\text{sty}, \text{sty})$

与知识库
 N 次相似度

How to Reduce Time Complexity?

核心思路：“层次过滤思想”



$O(N)$

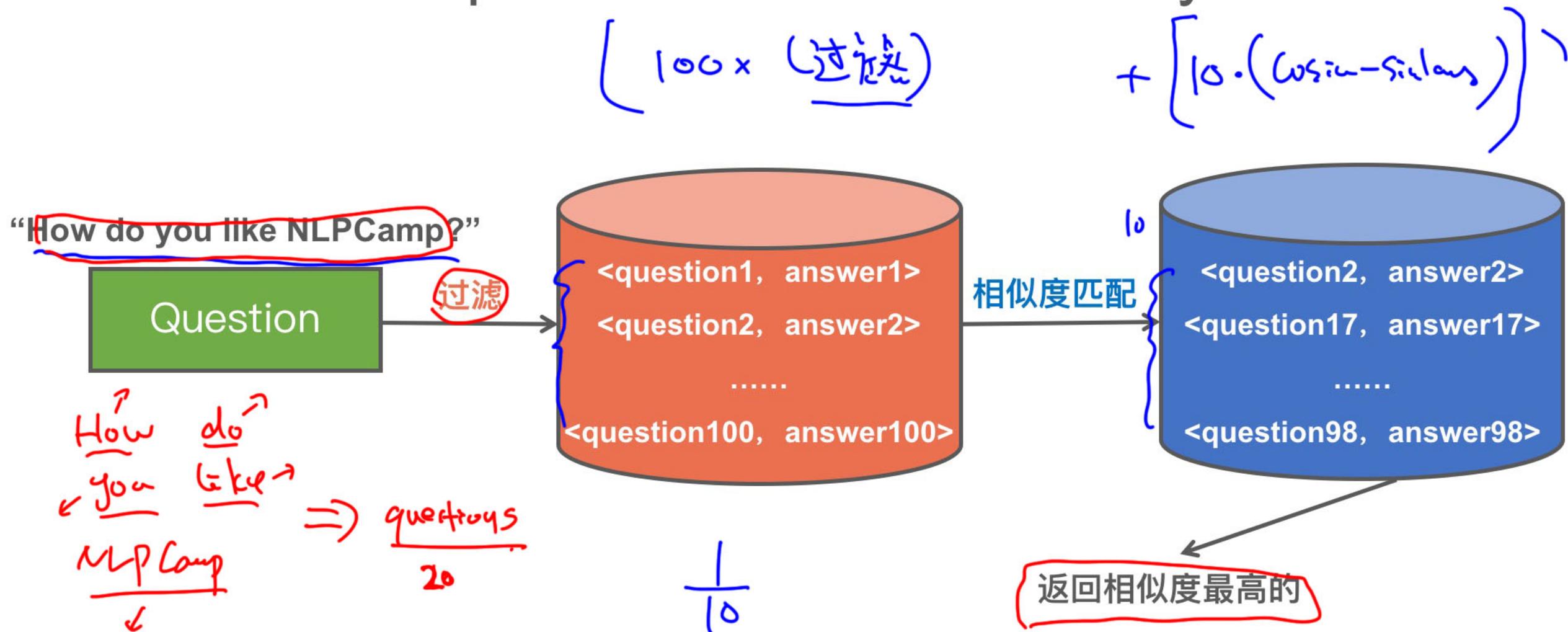
$N \neq \infty \Rightarrow O(N)$

实用性

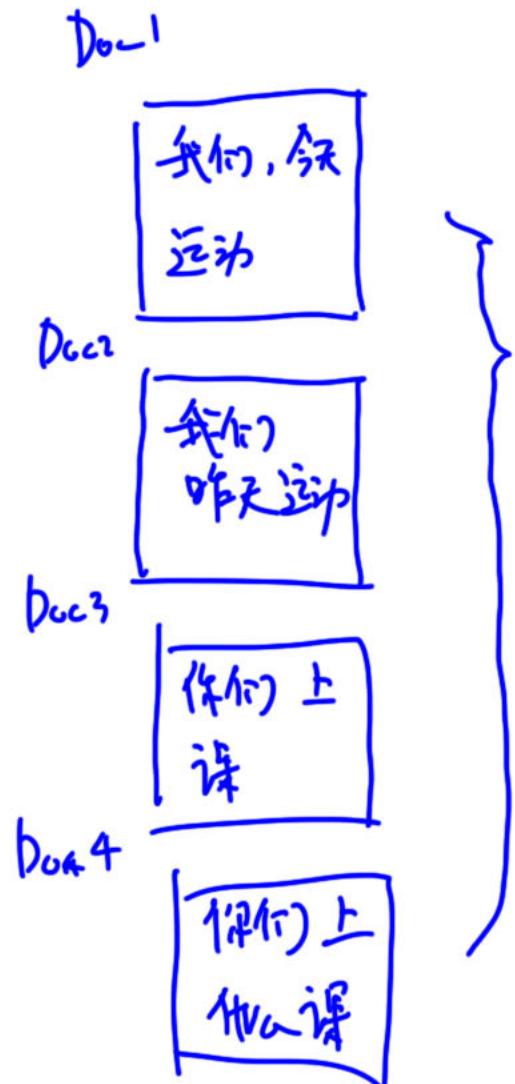
效率(过滤1) < 效率(过滤2)

< 效率(cosine-similarity)

Recap: Retrieval-based QA System

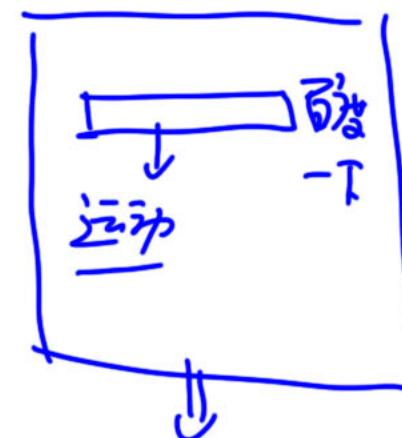


Introducing Inverted Index



词典: [我们, 今, 天, 运, 动, 昨, 天, 上, 课, 什, 么]

- 我们: [Doc1, Doc2]
- 今: [Doc1]
- 运: [Doc1, Doc2]
- 天: [Doc2]
- 上: [Doc3, Doc4]
- 课: [Doc3, Doc4]
- 什: [Doc4]



return Doc1 n Doc2

输入: 我们, 运动
我们: [Doc1, Doc2]
运动: [Doc3, Doc4]
我们 & 运动 = φ
[Doc1, Doc2, Doc3, Doc4]
Ranking