

Word2Vec

1.

什么是Word2Vec和Embeddings?

Word2Vec是从大量文本语料中以无监督的方式学习语义知识的一种模型，它被大量地用在自然语言处理（NLP）中。那么它是如何帮助我们做自然语言处理呢？Word2Vec其实就是通过学习文本来用词向量的方式表征词的语义信息，即通过一个嵌入空间使得语义上相似的单词在该空间内距离很近。Embedding其实就是一个映射，将单词从原先所属的空间映射到新的多维空间中，也就是把原先词所在空间嵌入到一个新的空间中去。

我们从直观角度上来理解一下，cat这个单词和kitten属于语义上很相近的词，而dog和kitten则不是那么相近，iphone这个单词和kitten的语义就差的更远了。通过对词汇表中单词进行这种数值表示方式的学习（也就是将单词转换为词向量），能够让我们基于这样的数值进行向量化的操作从而得到一些有趣的结论。比如说，如果我们对词向量kitten、cat以及dog执行这样的操作： $kitten - cat + dog$ ，那么最终得到的嵌入向量（embedded vector）将与puppy这个词向量十分相近。

首先，考虑如何表示一个word

英语大概有约1300万单词。词向量就是将单词映射到词空间的一个点。

(注意，词向量在英文里有两个可以互相替换使用的说法：*word embeddings*和*word vectors*)

为什么使用词向量？一个最直观的原因是，我们可以找到一个N维（N小于1300万）的空间，足够编码我们所有的单词。

每一个维度可以编码一些含义，例如语义空间可以编码时态、单复数和性别。

one-hot向量 (one-hot vector) :

one-hot向量就是利用一个 $R^{|V| \times 1}$ 向量来表示单词。

$|V|$ 是词汇表中单词的数量。

一个单词在英文词汇表中的索引位置是多少，那么相对应的那一行元素就是1，其他元素都是0。

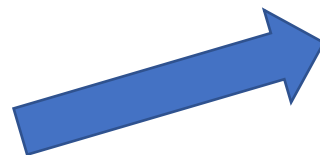
$$w^{Aardvark} = \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

$$w^{zebra} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix}$$



产生的问题是，任意两个单词的内积为0，无法区分他们之间具有的不同的相似性

$$(w^{hotel})^T w^{motel} = 0$$



$$v^c = \begin{bmatrix} 0.2 \\ 0.4 \\ 1.3 \\ \vdots \\ 0.9 \end{bmatrix}$$

那么如何得到这样的一个向量呢？

矩阵奇异值分解

DeepLearning

两种方法：SkipGram/CBow

以下面这个句子为例：{ "The" , "cat" , ' over" , "the' , "puddle" }。

通过这些上下文单词预测中间缺少的单词 "jump" , 这种模型我们就成为连续袋模型。

1. 将输入句子中的单词生成one-hot向量 (本文第二部分介绍过什么是one-hot向量) :

$$(x^{c-m}, \dots, x^{c-1}, x^{c+1}, \dots, x^{c+m} \in R^{|V|})$$

2. 通过将one-hot向量和输入矩阵 V 相乘, 得到输入单词的词向量:

$$\nu_{c-m} = Vx^{c-m}, \nu_{c-m+1} = Vx^{c-m+1}, \dots, \nu_{c+m} = Vx^{c+m}$$

3. 对上述词向量求平均得到: $\bar{\nu} = \frac{\nu_{c-m} + \nu_{c-m+1} + \dots + \nu_{c+m}}{2m} \in R^n$

4. 上面的平均词向量与输出矩阵进行点乘运算, 得到得分向量 (score vector)

MathJax maximum macro substitution count exceeded; is there a recursive macro c

, 我们知道, 两个向量约相似, 点乘得到的分数越高, 因此将会使相似的词互相靠近, 从而得到较高的分数。

5. 将分数转换成概率 $\bar{y} = softmax(z) \in R^{|V|}$ 。 ' softmax' 就是对 \bar{y} 做如下运算

$$\frac{e^{\bar{y}_i}}{\sum_{k=1}^{|V|} e^{\bar{y}_k}}, \text{ 使其每一个元素在}[0,1]\text{范围内, 且和为1。}$$

我们希望 \bar{y} 与真实的 y 匹配, 也就是真实的 y 的one-hot向量。

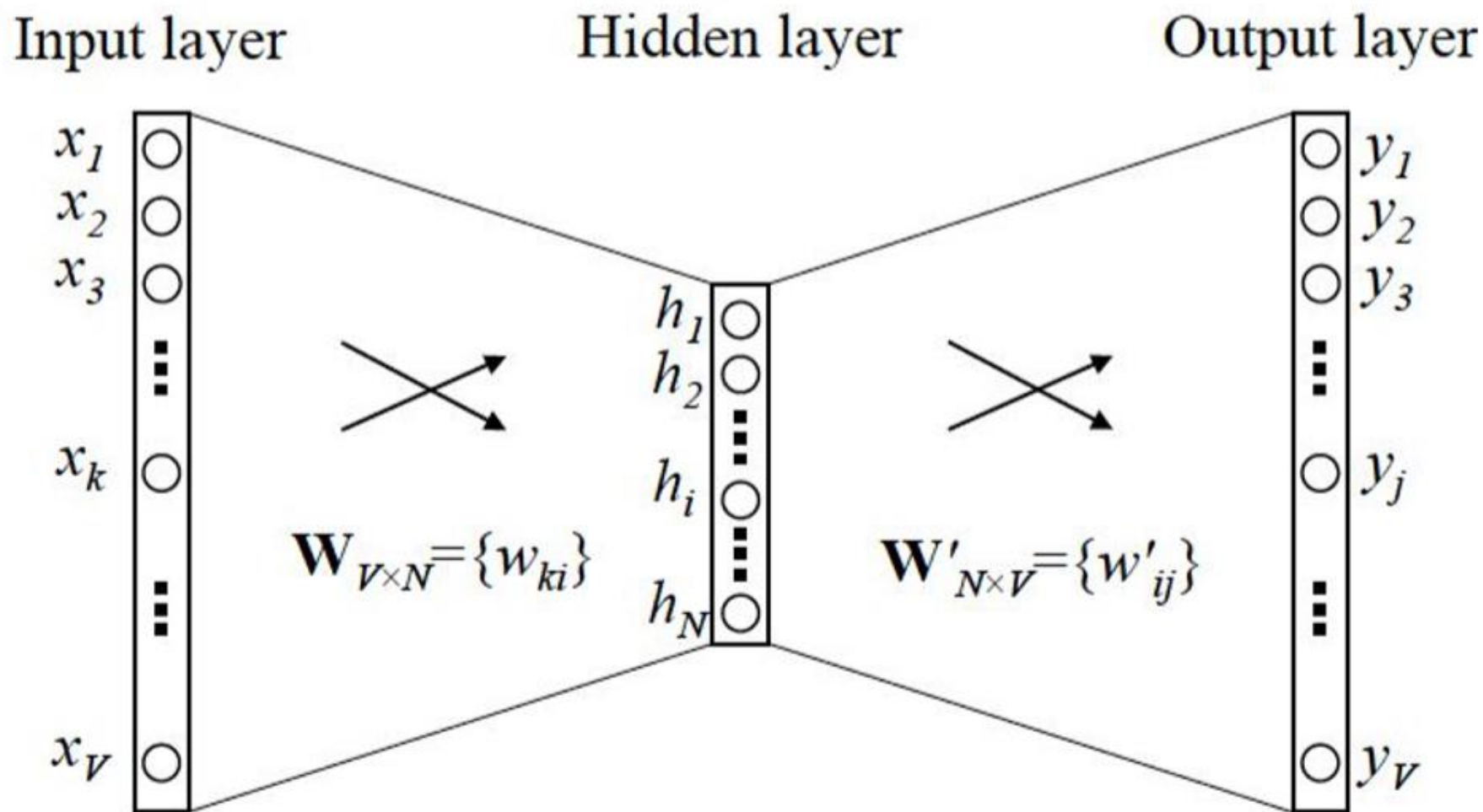
首先, 为了计算预测值的准确性, 我们先定义一个测量函数, 这里, 我们选择一个非常流行的测量函数: 交叉熵 (cross entropy)

$$H(\bar{y}, y) = - \sum_{j=1}^{|V|} y_j \log \bar{y}_j$$

Skip-Gram模型与连续袋模型正好相反。

连续袋模型是根据周围的词来预测中心的词，而Skip-Gram模型是根据中心的词预测周围的词。

Skip-gram Model



Demo的地址

https://github.com/Adoni/word2vec_pytorch/blob/master/model.py

```
word_embedding.txt - Visual Studio Code
word_embedding.txt x
1 8934 100
2 如何 0.095310114 -0.14539735 -0.09930189 0.72424227 0.39578062 -0.56750494 -0.25136185 -0.41487795 -0.13897215 -0.025657268 0.58856744 0.20525631 -0.21127132
3 装逼 0.12051985 -0.102363154 -0.008421278 0.15624484 0.7576986 0.19949712 -0.6565868 -0.15037183 -0.08343341 0.26319775 0.3988422 -0.08409254 0.25988683 0.14
4 装 -0.04010059 -0.20331523 0.16845343 0.3282443 0.50755745 -0.015349344 -0.45837075 0.094665214 -0.13542148 0.42710677 0.37323028 0.37945718 0.28138542 0.14
5 得 0.14839926 -0.44357243 0.11422428 -0.1352943 0.4727048 -0.21937385 -0.2139999 0.23827878 -0.5578143 0.35479656 0.2569447 0.04590395 0.23896088 0.07107336
6 别人 0.22392417 -0.32456285 0.105446324 0.078167304 0.27501112 -0.0038775087 -0.6373443 -0.28652796 -0.27973026 0.06367408 -0.069404624 0.29041064 0.483086 0
7 一 0.27620238 0.17761037 0.30916548 0.070706874 -0.3938096 -0.055420086 -0.79651594 0.019092567 0.22868624 0.09661173 0.23614109 0.4720921 0.3891802 0.19166
8 的 0.112183474 0.1450165 -0.051926028 0.123593025 0.3286875 -0.052716043 -0.15170953 0.02135227 0.06660297 -0.04773142 0.2804304 0.4639615 0.68142515 0.1214
9 ？ 0.61203104 -0.045311943 -0.27664408 -0.019538121 0.25580665 0.10690243 -0.39192888 0.08835911 0.25091937 -0.0013977016 0.48325312 0.03520442 0.19629812 -0
10 你 0.08241192 0.21917224 -0.06616307 0.003936575 0.6645792 0.43512845 -0.2296878 0.020709572 -0.5836536 -0.25745904 0.15010591 -0.02846796 1.0271162 0.35766
11 拥有 0.25903895 -0.19976307 -0.06376426 0.10867175 0.27492356 0.34095615 -0.17216538 -0.18372557 0.03699752 -0.1990261 0.042120337 0.3738848 0.09034352 0.23
12 哪些 0.1596949 -0.48724672 -0.2622621 0.6115895 0.78305393 0.07821899 -0.41349924 0.13206509 -0.17348069 0.87447244 1.1508789 0.77845275 0.7753424 0.15339428
13 为什么 -0.0032535086 0.3957269 -0.035097823 -0.20033433 0.1484299 0.23379499 -0.26663682 0.6125705 -0.5750434 0.4305 0.5440889 -0.06735996 0.2029098 0.429743
14 接吻 0.052189056 -0.020449338 -0.07212644 0.005872586 0.3443337 0.049990475 -0.4616933 -0.07608312 -0.048890736 0.052053574 0.17588001 0.30190825 0.16617176
15 时候 -0.24034216 -0.38797712 -0.09306655 0.30747157 0.6327519 0.3261789 -0.4557257 0.08258574 -0.42369413 -0.050618675 0.4087699 0.2724262 0.21528952 -0.0772
16 男生 0.19576955 -0.3140485 -0.11978105 -0.15464716 0.62988806 0.121418655 -0.2525586 -0.14947248 -0.11527305 -0.31995088 0.07593282 0.48156917 0.21842492 0.2
17 女生 0.28566796 -0.23167193 -0.12660341 -0.1936253 0.57374805 0.047900192 -0.2077223 -0.06689489 -0.10753916 -0.21258056 0.12636657 0.28805327 0.1654635 -0.6
18 胸部 0.05456499 -0.193639 -0.06022215 -0.068425536 0.37414482 -0.017197428 -0.46060854 -0.1795102 -0.006473076 -0.0089789685 0.10294023 0.2585926 0.0924423 0
19 汽车 0.03839941 -0.1348901 -0.27590948 0.26693073 0.10294195 -0.00035869886 -0.3902156 0.24938342 0.03744077 0.12218074 0.55123353 0.33543172 0.5700409 -0.1
20 行业 0.16517642 0.05251761 -0.46800625 0.24080999 0.2780133 0.041723933 0.5565421 -0.22911958 0.05125378 -0.057586923 0.5280789 0.15964231 0.6337026 -0.4736
21 有 0.25330573 -0.39740396 -0.5136794 -0.16841248 0.13030499 0.026022831 -0.184139 -0.539832 0.24745719 0.24414149 0.8348171 0.22523405 0.21827206 0.19186752
22 笑话 0.20197171 -0.11150992 -0.06669846 0.11104404 0.28363115 0.11517108 -0.54979926 0.042639703 -0.16402404 0.1683029 0.44997212 0.40166914 0.45023468 -0.00
23 快速 0.08547922 -0.06389757 -0.24138823 0.19590281 0.4034331 0.105244495 -0.15352026 -0.11542776 0.057633348 -0.03934344 0.3048158 0.43123242 0.34863147 -0.3
24 去除 -0.01728689 -0.091094375 -0.07842307 0.11838398 0.18317823 -0.013177681 -0.40312982 -0.09661373 -0.1513989 0.016840156 0.21665174 0.23247541 0.24302368
25 前 0.17596051 0.163249 0.01849693 0.34114102 0.20873912 -0.19223914 -0.43874812 -0.0035307084 -0.1294078 0.2125397 -0.038686793 0.22510204 0.0966728 -0.0207
26 玻璃 0.02809285 0.04765922 -0.042685654 0.18886809 0.15371433 0.036027964 -0.43408158 -0.080201164 -0.16819742 0.04143994 0.17863162 0.3192556 0.2244503 0.00
27 上 0.48287383 0.06644198 -0.0051682782 0.45631358 0.6744782 0.66312116 -0.4998817 0.4851277 -0.27309367 0.39711517 0.22631219 -0.13220626 0.10027322 -0.2030
28 冰 -0.06295095 0.15620072 -0.03945921 0.14534976 -0.046646666 -0.18259606 -0.43469313 -0.084874906 -0.09541632 0.106739715 0.63312197 0.042122126 0.18570083
29 发动机 0.06608668 0.015852084 -0.34292862 0.34148848 0.10779312 -0.09877443 -0.3716562 0.050943356 -0.10632809 0.024954205 0.2982568 0.1269363 0.4537813 -0.0
30 漏 0.02820562 -0.065083005 -0.054954864 0.07529259 0.17977105 0.03804293 -0.3705521 -0.05114434 -0.11862666 0.043307357 0.2311383 0.2886472 0.3185508 -0.004
31 油 0.09525751 -0.079573154 -0.021433797 0.059250582 0.3294751 0.06882176 -0.48361441 0.005241505 0.057705387 0.029175946 0.2201137 0.43998632 0.3754099 -0.08
32 的 0.021990178 -0.15861888 -0.08671903 0.18508191 0.2865082 0.06852269 0.086880885 -0.10977919 -0.004417409 0.003241688 0.45452163 0.48846826 0.3866135 0.034
33 问题 -0.15140684 0.4009903 -0.4093234 0.3732281 0.18605483 -0.22999936 -0.47437888 -0.24694358 0.0024290164 0.15068954 0.121026106 -0.268734 0.86847705 -0.36
34 大 0.099050894 0.039894722 -0.6046009 -0.026246954 0.30907026 0.2252479 -0.85428435 -0.6912352 -0.17802347 -0.1514285 0.42167008 0.22830467 0.60730296 -0.08
35 吗 0.24396285 0.50626624 -0.6836878 0.39321658 0.15681495 0.15258609 -0.48499066 -0.065966725 0.053084705 0.110154785 0.34262633 0.47591397 0.0036833996 0.0
36 在 0.09509084 -0.446293 -0.30237514 0.5174772 0.54933417 0.82545996 -0.29103243 0.35713178 0.34397087 0.23922448 0.8471535 0.14824255 0.41070172 -0.15874638
37 达到 -0.13745774 -0.26810902 -0.2600764 0.511728 0.20369317 -0.05254516 -0.103267185 0.23568349 0.04897126 0.06344485 0.07505992 0.44279662 0.38214853 0.0393
38 极限 0.016238617 0.020571552 -0.19256672 0.19612874 0.10353231 0.095191 -0.21817622 -0.041644637 -0.03048089 -0.070371665 0.29137024 0.25024432 0.3640875 -0.
39 之前 -0.064108625 -0.29179123 0.010118262 0.33190507 0.28938958 0.35458753 -0.5685632 0.11780529 -0.22777288 -0.25405028 0.3640842 0.39366376 0.5004031 0.120
40 或 0.43980885 0.28814092 -0.5022895 0.05341447 0.25642416 0.3743128 -0.6540276 -0.45082676 0.22251064 0.18720852 0.40923035 0.3985554 0.3497653 -0.18965007
41 车辆 -0.054185912 -0.031308692 -0.10310428 0.1763499 0.25264448 -0.07866875 -0.31379417 -0.028955752 0.18938614 -0.14837427 0.33348647 0.35743383 0.25720742
```

负采样

<https://www.cnblogs.com/pinard/p/7249903.html>

Google 官方的word2Vec项目

<https://github.com/tmikolov/word2vec>