

Text summarization part 1

文本摘要应用场景：QA，论文摘要，上市公司的报告，^(情)舆论控制，Newsletter

algorithm engineer: 算法基础，代码功底，业务理解能力和应用能力

良好的产品sense和业务owner意识

Language

Denotational Semantics meaning

Wordnet: ① Great but missing nuance ② Missing new meaning of word
③ Subjective ④ requires human labor to create and adapt
⑤ Can't compute accelerate word similarity

Representing words as discrete symbols:

ex: one-hot encode



problems: ① 词表大且增长快 ② 很难表示相近的含义：motel, hotel
③ 不具有通用性

用上下文表达它们的含义

N-gram:
N=1 unigram
N=2 bigram
N=3 trigram

N-gram LM

out of vocabulary

problem of N-gram: 数据稀疏，难免会呈现 OOV 问题

随着 n 的增大，参数空间呈指数增长（维度灾难）

缺少长期依赖，只能建模到前 n-1 个词

无法表示一词多义（语义鸿沟）

word vector \Rightarrow also called word embeddings or (neural) word representations

Semantic and syntactic relationships ex: king & Queen 他们的距离是很接近的

Neural Probabilistic Language Model (2003)

Problem of NPL model: 有限的前文信息
计算量过大 词向量是副产品
 \downarrow
在 softmax 和 tanh 中间全连接层
计算量很大

Word2Vec:

$$\text{Likelihood } L(\theta) = \prod_{t=1}^T \prod_{j=0}^{m-1} P(W_{t+j} | W_t, \theta)$$

$$\text{objective function } J(\theta) = -\frac{1}{T} \log L(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{j=0}^{m-1} \log P(W_{t+j} | W_t, \theta)$$

minimize objective function \hookrightarrow Maximizing predictive accuracy

how to calculate $P(W_{t+j} | W_t, \theta)$?

Answer: use two vectors per word w_i :

v_w when w is a center word

u_w when w is a context word

$$P(w|c) = \frac{\exp(u_w^T v_c)}{\sum_{v \in V} \exp(u_w^T v)}$$

Dot product compares similarity of u and v , $u^T v = u \cdot v = \sum_i u_i v_i$
Larger dot product = larger probability
Normalize over entire vocabulary to give probability distribution

ex: softmax function $R^n \rightarrow (0, 1)^n$

$$\text{softmax}(x_i) = \frac{\exp(x_i)}{\sum_{j=1}^n \exp(x_j)} = p_i$$

"max" because amplifies probability of largest x_i → return a distribution

"soft" because still assigns some probability to smaller x_i

frequent used in DL

Word2Vec:

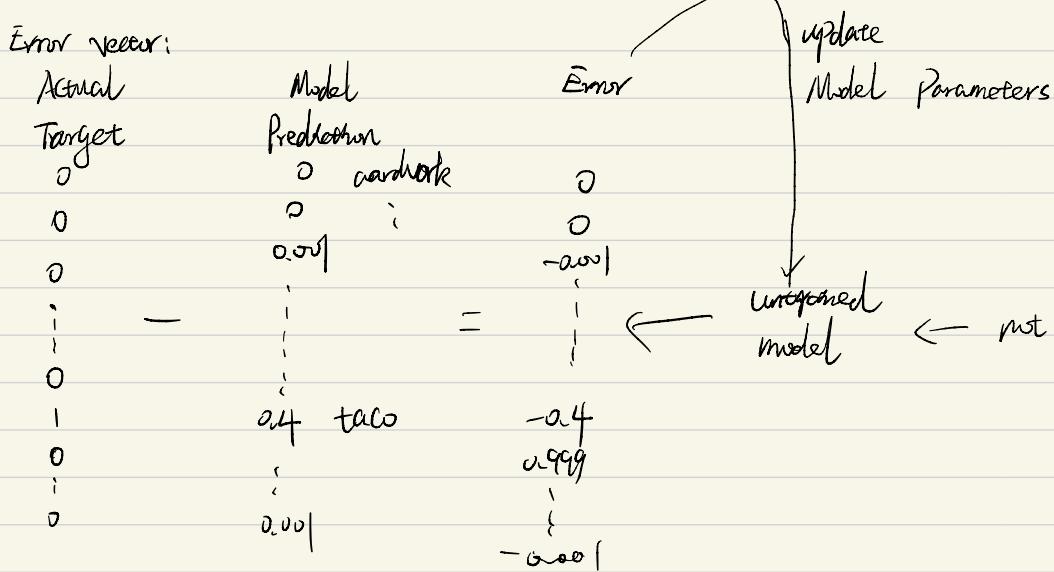
① CBOW 多个输入固定字来预测一个中心词

② skip-gram 一个中心词来预测周边的单词

Extract first sample

not → Untrained Model

Task:
Predict neighboring word



Diff between skip-Gram & CBOW:

CBOW is simpler & faster; skip-gram 学习到的语义更多
CBOW 更倾向学习同义词的语义

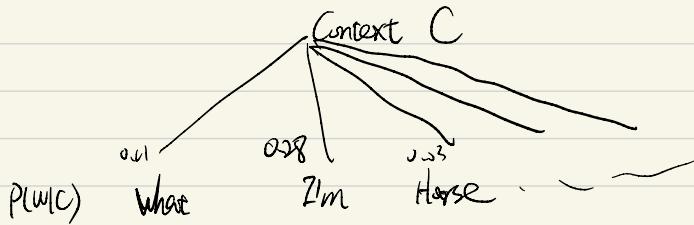
Problem of word2vec:

类别太多，计算量大



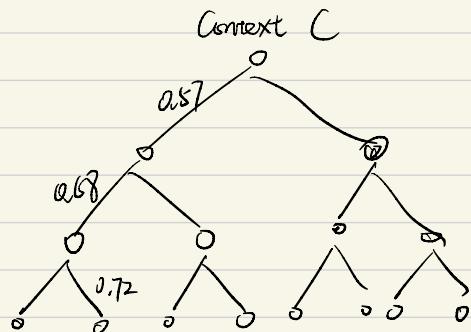
解决方案：Hierarchical Softmax (H-Ssoftmax)

Switch Softmax to Binary Tree



You must compute all of the terminal leaves

$O(n)$ to $O(\log(n))$



$$P(I'm|C) = 0.57 \cdot 0.68 = 0.72 = 0.28$$

$P(w|C)$ what I'm Horse

Encode Node

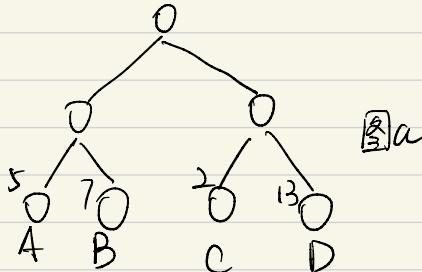
Construct Tree & Huffman Encoding

- Random construct

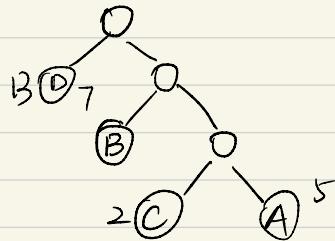


Term Frequency Huffman Tree

Advantage for Huffman Encoding



图a



图b

带权路径长度

$$\text{图a: } WPL = 5*2 + 7*2 + 2*2 + 13*2 = 54$$

$$\text{图b: } WPL = 5*3 + 2*3 + 7*2 + 13*1 = 48$$

这类方法都是 approximating the softmax，大约提高了50X的速度

Problem of Hierarchical Softmax:

当预测的词是生僻词时(频率使用比较低的时候)走的路径很长，计算量非常大

Negative sampling(负采样):

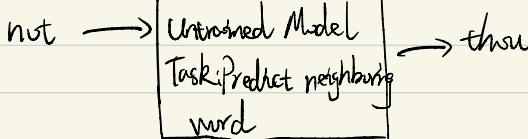
在实际场景中，embedding一般是已经训练好了的

- Softmax : 1) look up embedding
- 2) calculate prediction
- 3) project to output vocabulary

Switch to logistic regression Task:

一开始是利用上下文来预测：

ex:



To:

not →

then →

Untrained Model

Task:

Are the two words
neighbors?

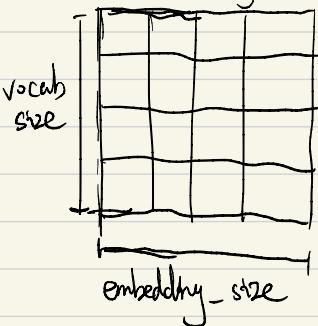
→ 0.90 (预测正确概率, 最好接近于1)

怎么进行负采样: 从 vocabulary 中而不是上下文的单词中拿出单词,
构建一个负样本。

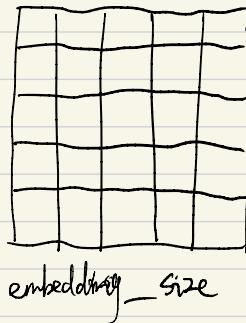
skip-Gram with negative Sampling (SGNS):

start with two Random Matrices

Embedding



Context



target - sigmoid)

input word	output word	target	input + output	sigmoid()	Error()
not	throw	1	0.2	0.55	0.45
not	acron	0	-1.11	0.25	-0.25
not	taco	0	0.74	0.68	-0.68

超参数 Hyperparameters Gensim default 为 5 : window size

smaller window sizes (2-15) lead to embeddings where high similarity scores between two embeddings indicates that the words are interchangeable

Large window sizes (15-50, over more) lead to embeddings where similarity is more indicative of relatedness of the words.



实际情况是可以通过 word tree 来看 embedding 的关系

Negative samples: Gensim 中 default 为 5

正常情况下，语料或数据集越大，negative sampling 算法越小；反之亦然