# NLP 基础算法梳理(One-hot, Word2vec, Elmo, Bert, GPT)
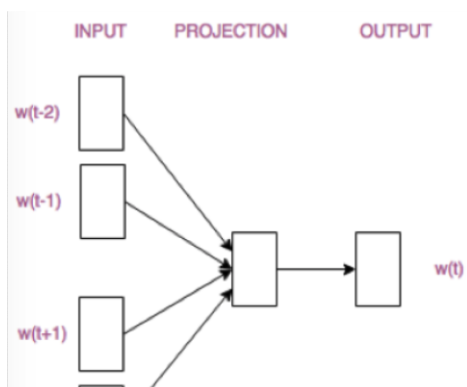
Friday, May 28, 2021    2:13 AM

教程来自Hongyi Li的教程和Shusen Wang的教程以及Jay Alammar的博客

自从seq2seq模型引入到NLP领域之后，后续的模型基本上都延续了这encoder-decoder的结构。

其中 encoder 的目的是将字词进行编码，得到字词在多维向量空间representation，这个过程叫word embedding

## 1. Embedding的方法更迭

- One-hot encoding：0,1对词库里的所有词进行编码。
  - 缺点是词的数量增多的时候会出现维度的爆炸；此外还有矩阵稀疏的问题。

- Word2vec，推进了NLP领域的发展
  - CBOW



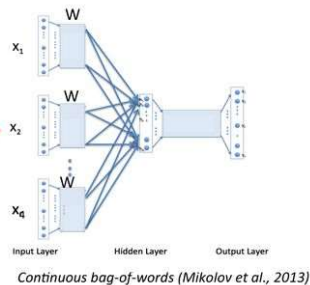An example of CBOW Model



Corpus = { I drink coffee everyday }

Initialize:
$$W = \begin{bmatrix} 1 & 2 & 3 & 0 \\ 1 & 2 & 1 & 2 \\ -1 & 1 & 1 & 1 \end{bmatrix}$$

Ex:
$W^{drink} = [0,1,0,0]$

$$\begin{bmatrix} 1 & 2 & 3 & 0 \\ 1 & 2 & 1 & 2 \\ -1 & 1 & 1 & 1 \end{bmatrix}\begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \quad Wx_2 = v_2 \quad = \begin{bmatrix} 2 \\ 2 \\ 1 \end{bmatrix}$$
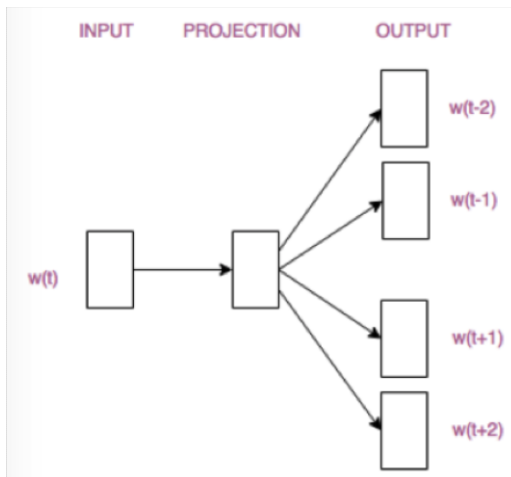
Continuous bag-of-words (Mikolov et al., 2013)

优点：向量运算有实际的意义，如King-Man+woman = queen
缺点：每个单词只有一个向量来表征，没法分辨同义词

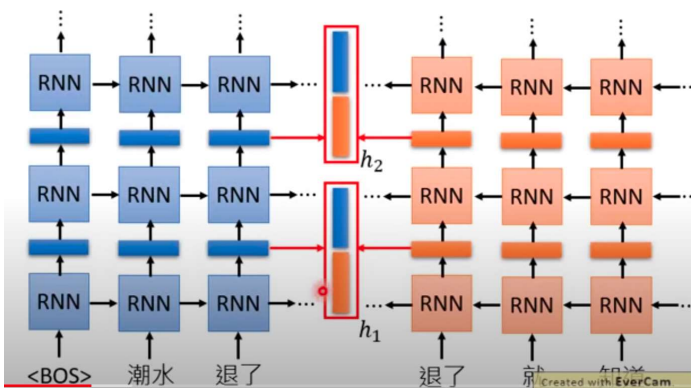为了减少训练的时间，有同义词之间只是微调参数，引入了霍夫曼树的方法

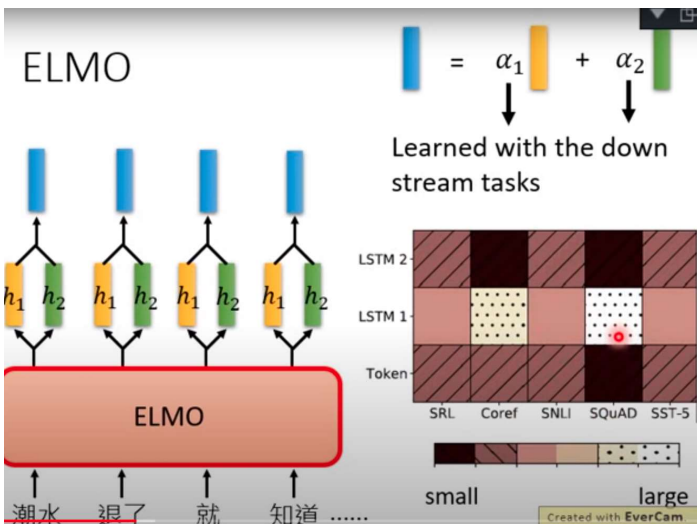  - Skip-gram

**优点：解决了contextual的问题，根据上下文意思的不同，每个词的Embedding是不一样的**

**缺点：LSTM是seq2seq的，速度慢**

- ELMo（Embedding from language models）



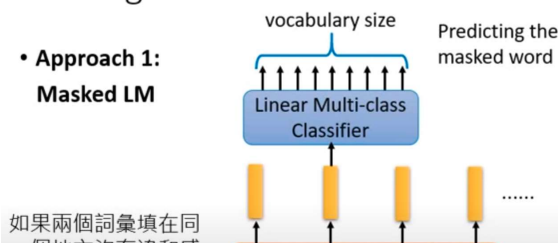每一层的LSTM能够生成一个latent representation

不同任务的a1和a2是不一样的，这两个参数根据不同的任务来训练

右下角的图表示的是哪一层的参数对任务的影响最大



ELMO

$$\blacksquare = \alpha_1 \blacksquare + \alpha_2 \blacksquare$$

Learned with the down stream tasks

- BERT（Bidirectional Encoder representation from transformers）

两种训练方法，一种是Masked LM，预测句子中被遮盖住的某 **一起训练的**

一般情况下两种都是



**Approach 2: Next Sentence Prediction**

No

Linear Binary Classifier

[SEP]: the boundary of two sentences

[CLS]: the position that outputs classification results

Approaches 1 and 2 are used at the same time.

BERT

[CLS] 醒醒 吧 [SEP] 眼睛 業障 重

Training of BERT

vocabulary size

Predicting the masked word

- **Approach 1: Masked LM**

Linear Multi-class Classifier

如果两个词彙填在同一個地方沒有違和感

## Training of BERT

- **Approach 1: Masked LM**

vocabulary size — Predicting the masked word

Linear Multi-class Classifier

如果兩個詞彙填在同一個地方沒有違和感

那它們就有類似的 embedding

BERT

[CLS] 醒醒 吧 [SEP] 眼睛 業障 重

## How to use BERT – Case 1

B

class

Linear Classifier → Trained from Scratch

BERT → Fine-tune

[CLS] $w_1$ $w_2$ $w_3$

sentence

Input: single sentence, output: class

Example: Sentiment analysis (our HW), Document Classification

## How to use BERT – Case 2

class  class  class

Linear Cls  Linear Cls  Linear Cls

BERT

[CLS] $w_1$ $w_2$ $w_3$

sentence

Input: single sentence, output: class of each word

Example: Slot filling

arrive  Taipei  on  November  2nd

other  dest  other  time  time

## How to use BERT – Case 3

Class

Linear Classifier

BERT

[CLS] $w_1$ $w_2$ [SEP] $w_3$ $w_4$ $w_5$

Sentence 1    Sentence 2

Input: two sentences, output: class
Example: Natural Language Inference
Given a "premise", determining whether a "hypothesis" is T/F/ unknown.

## How to use BERT – Case 4

Learned from scratch

$s = 2$  $e = 3$
The answer is "$d_2 d_3$".

0.1   0.2   0.7

Softmax

dot product

BERT

[CLS] $q_1$ $q_2$ [SEP] $d_1$ $d_2$ $d_3$

question    document

## What does BERT learn?

https://arxiv.org/abs/1905.05950
https://openreview.net/pdf?id=SJzSgnRcKX

POS       $K(\Delta) = 1.60$    $K(s) = 0.19$
Consts.   $K(\Delta) = 1.57$    $K(s) = 0.83$
Deps.     $K(\Delta) = 1.15$    $K(s) = 0.87$
Entities  $K(\Delta) = 1.61$    $K(s) = 0.06$

F1 Score    Expected layer & center of gravity

这个分层，就是把每一层encoder的hide state提取出来

其实各个层级的hiden state加起来显示的F1是不一样的

FROM PAPER:
Dissecting Contextual Word Embeddings: Architecture and Representation
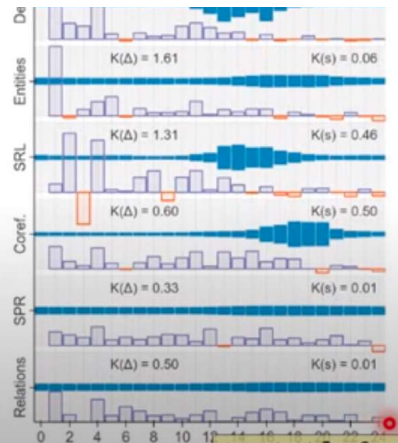
WHAT DO YOU LEARN FROM CONTEXT, PROBING FOR

BERT Rediscovers the Classical NLP Pipeline

Evaluation as word representations
- MultiNLI dataset
- Semantic Role Labeling, Ontonotes dataset
- Constituency parsing, Penn Treebank

| | F1 Scores | | Expected layer & center-of-gravity |
|---|---|---|---|
| | ℓ=0 | ℓ=24 | 0  2  4  6  8  10  12  14  16 |
| POS | 88.5 | 96.7 | 3.39 — 11.68 |
| Consts. | 73.6 | 87.0 | 3.79 — 13.06 |
| Deps. | 85.6 | 95.5 | 5.69 — 13.75 |
| Entities | 90.6 | 96.1 | 4.64 — 13.16 |
| SRL | 81.3 | 91.4 | 6.54 — 13.63 |
| Coref. | 80.5 | 91.9 | 9.47 — 15.80 |
| SPR | 77.7 | 83.7 | 9.93 — 12.72 |
| Relations | 60.7 | 84.2 | 9.40 — 12.83 |

Entities — K(Δ) = 1.61    K(s) = 0.06
SRL — K(Δ) = 1.31    K(s) = 0.46
Coref. — K(Δ) = 0.60    K(s) = 0.50
SPR — K(Δ) = 0.33    K(s) = 0.01
Relations — K(Δ) = 0.50    K(s) = 0.01

Evaluation as word representations
- MultiNLI dataset
- Semantic Role Labeling, Ontonotes dataset
- Constituency parsing, Penn Treebank
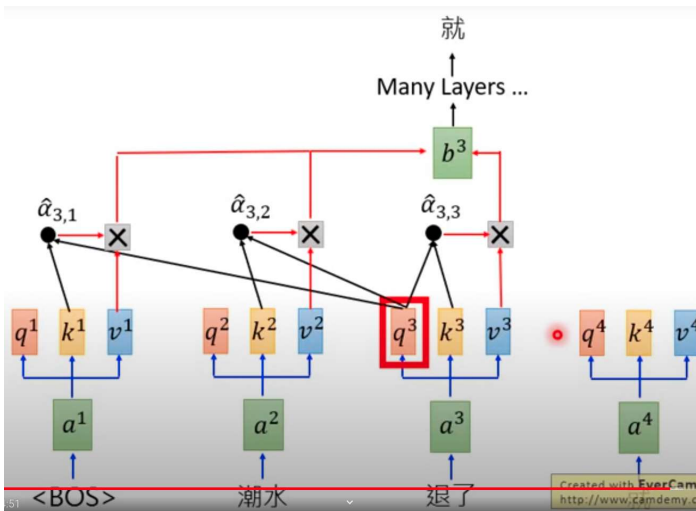- Name entity recongnition: NER

Contextual similarity
- Intra-sentence similarity
- Span representations
- Unsupervised pronominal coref

Probing contextual information
- POS：part-of-speech tagging
- Constituency parsing

- GPT（Generative Pre-training）是Transformer的decoder



- **Reading Comprehension**

$$d_1, d_2, \cdots, d_N, "Q:", q_1, q_2, \cdots, q_M, "A:"$$

- **Summarization**   $d_1, d_2, \cdots, d_N, "TL;DR:"$

- **Translation**

| English sentence 1 | = | French sentence 1 |
|---|---|---|
| English sentence 2 | = | French sentence 2 |
| English sentence 3 | = | |

## Multi-head Self-attention — Different types of relevance

$$q^{i,1} = W^{q,1}q^i$$
$$q^{i,2} = W^{q,2}q^i$$



$$q^i = W^q a^i \qquad \text{(2 heads as example)}$$

2. Attention

- (NNLM) A Neural Probabilistic Language Model
- (Word2Vec) Distributed Representations of Words and Phrases and their Compositionality
- (GloVe) GloVe: Global Vectors for Word Representation
- (ELMo) Deep contextualized word representations
- (GPT) Improving Language Understanding by Generative Pre-Training
- (BERT) BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding