Week1.    CS224n : NLP with Deeping Learning

Lecture 1 : Introduction and Word Vectors

In traditional NLP, words as discrete symbols : <u>a localist representation</u>

Means one 1s, the rest 0s      one-hot vectors

vector dimension is big (500,000+)

But also no <u>similarity</u> for one-hot vectors

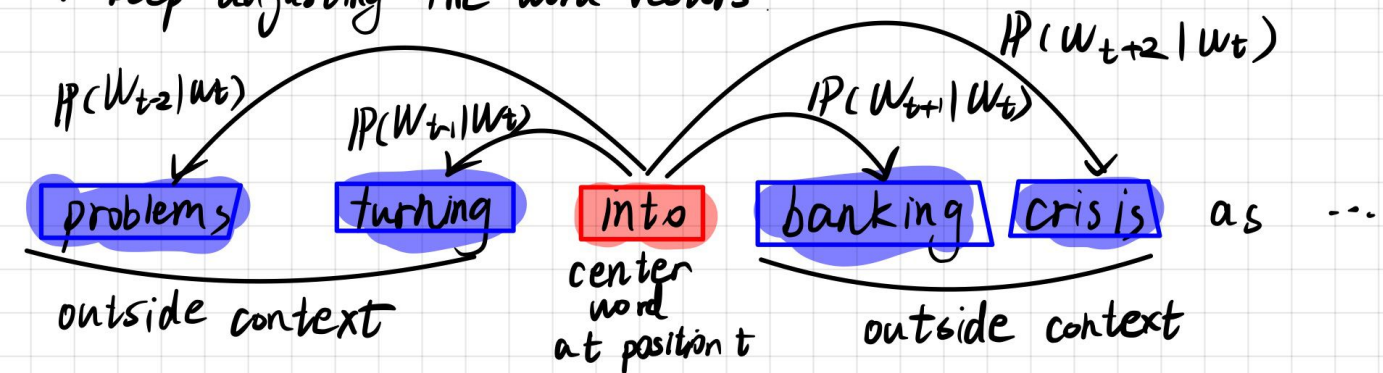**\* learn to encode similarity into the vectors themselves**

Representing words by their context.

"A word's meaning is given by the words that frequently appear near-by"

Word Vectors : a dense distributed vector

# Word2vec:

- a large corpus
- each word by a vector
- has center word "c" and context outside word "o"
- use similarity of the word vectors for "c" and "o" to calculate the probability of "o" given "c".
- keep adjusting the word vectors



$P(W_{t+2}|W_t)$

$P(W_{t-2}|W_t)$    $P(W_{t-1}|W_t)$      $P(W_{t+1}|W_t)$

problems    turning    Into    banking    crisis    as    ...

outside context      center word at position t      outside context

For each position $t = 1, \dots, T$, predict contexts words within a window of fixed size $m$, given center word $W_j$.

Likelihood $L(\theta) = \prod\limits_{t=1}^{T} \prod\limits_{\substack{-m \le j \le m \\ j \ne 0}} P(W_{t+j}|W_t ; \theta)$

$\theta$ is all parameters to be optimized

The objective function $J(\theta)$ is the average Log likelihood:

$$J(\theta) = -\frac{1}{T} L(\theta) = -\frac{1}{T} \sum\limits_{t=1}^{T} \sum\limits_{\substack{-m \le j \le m \\ j \ne 0}} \log P(W_{t+j}|W_t ; \theta)$$

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

Answer: 2 vectors • $V_w$ when $w$ is a center word
• $U_w$ when $w$ is a context word
Then for a center word $c$ and a context word $o$ :

$$P(o|c) = \frac{\exp(U_o^T V_c)}{\sum_{w \in V} \exp(U_w^T V_c)}$$

Dot product produces similarity

sum to normalize

• This is an example of Softmax function $\mathbb{R}^n \mapsto \mathbb{R}^n$

$$softmax(X_i) = \frac{\exp(X_i)}{\sum_{j=1}^n \exp(X_j)} = P_i$$

• "max" means amplifies probability of largest $X_i$
• "soft" still assign some probability to smaller $X_i$
• Frequent used in Deep Learning

To train the model, Compute all $\boxed{2d}$ vector gradients.

(recall every word has 2 vectors)

**Blackboard Derivation**

$$\max \quad L(\theta) = \prod_{t=1}^{} \prod_{\substack{-m \le j \le m \\ j \ne 0}} P(W_{t+j} | W_t; \theta)$$

$$\min \quad J(\theta) = -\frac{1}{T} \sum_{t=1}^{T} \sum_{\substack{-m \le j \le m \\ j \ne 0}} \log P(W_{t+j} | W_t; \theta)$$

$$P(o|c) = \frac{\exp(U_o^T V_c)}{\sum_{w \in V} \exp(U_w^T V_c)}$$

$\langle$ useful tips : $\frac{\partial X^T a}{\partial X} = \frac{\partial a^T X}{\partial X} = a \rangle$

$$\frac{\partial}{\partial V_c} \log \frac{\exp(U_o^T V_c)}{\sum_{w \in V} \exp(U_w^T V_c)} = \frac{\partial}{\partial V_c} \log (\exp(U_o^T V_c)) - \frac{\partial}{\partial V_c} \log \sum_{w \in V} \exp(U_w^T V_c)$$

$$= \frac{\partial}{\partial V_c} U_o^T V_c - \frac{\frac{\partial}{\partial V_c} \sum_{w \in V} \exp(U_w^T V_c)}{\sum_{w \in V} \exp(U_w^T V_c)}$$

$$= U_o - \frac{\sum_{x=1}^{V} \exp(U_x^T V_c) \cdot U_x}{\sum_{w \in V} \exp(U_w^T V_c)}$$

$$= U_o - \sum_{x=1}^{V} P(x|c) U_x$$

Current context word        expected context word

home derivation

$$\frac{\partial}{\partial u_o} \log \frac{\exp(u_o^T v_c)}{\sum_{w \in V} \exp(u_w^T v_c)} = \frac{\partial}{\partial u_o} u_o^T v_c - \frac{\partial}{\partial u_o} \log \sum_{w \in V} \exp(u_w^T v_c)$$

$$= v_c - \frac{\frac{\partial}{\partial u_o} \sum_{w \in V} \exp(u_w^T v_c)}{\sum_{w \in V} \exp(u_w^T v_c)}$$

$$= v_c - \frac{\frac{\partial}{\partial u_o} \exp(u_o^T v_c)}{\sum_{w \in V} \exp(u_w^T v_c)}$$

$$= v_c - \frac{\exp(u_o^T v_c)}{\sum_{w \in V} \exp(u_w^T v_c)} v_c$$

$$= v_c - P(o \mid c) \cdot v_c$$

Lecture 2: Word Vectors and Word Senses

$* \nabla_\theta J(\theta)$ very sparse in SGD
- only update the word vectors that actually appear.

Two model variants:

1. Skip-grams (SG)
   predict context ("outside") words (position independent) given center words
2. Continuous Bag of Words [CBOW]
   predict center word from (bag of) context words

Additional efficiency in training: Negative sampling

In hw 2: the skip-gram model with negative sampling

window-based coourrence matrix
- I like deep learning
- I like NLP
- I enjoy flying

- increase in size
- sparsity
- much storage

| | I | like | enj | dep | leariy | NLP | flyiŋ | . |
|---|---|---|---|---|---|---|---|---|
| I | 0 | 2 | 1 | 0 | 0 | 0 | 0 | 0 |
| like | 2 | 0 | 0 | 1 | 0 | 1 | 0 | 0 |
| enj | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 |
| dep | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 |
| learniŋ | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 |
| NLP | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 |
| fly | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 |
| . | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 0 |

Solution: Low dimensional vectors

Method: SVD (hw1)    Factor X into $\underline{U \Sigma V^T}$

Method: Hacks to X            · $\min(X, t)$, $t \approx 100$
                              · ramped windows, count closer words more

Count-based v.s. direct prediction

| |
|---|
| · LSA, HAL |
| · COALS, Hellinger-PCA |
| · Fast training |
| · Efficient use of statistic |
| · primarily used to capture word similarity |
| · Disproportionate importance given to large corpus |

| |
|---|
| · skip-gram / CBoW |
| · NNLM, HLBL, RNN |
| · Scale with corpus size |
| · Inefficient use of statistics |
| · Generate improved performance |
| · capture complex patterns |

Evaluation in NLP: intrinsic v.s. extrinsic

Word Sense ambiguity:

- common words

- words that exist for a long time