

Word2Vec

20424008 - Dương Mạnh Cường



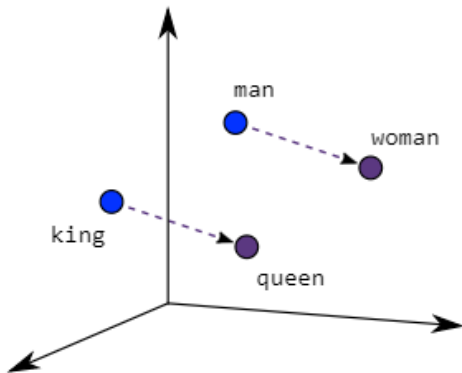
*Đại học Khoa Học Tự Nhiên
ĐHQG Thành phố Hồ Chí Minh*

14/02/2022

Nội dung

- 1 Word2Vec model
 - CBOW model
 - CBOW model với một context word
 - CBOW model với multiple context words
 - Skip-gram model
- 2 Các chiến lược cải thiện hiệu suất
 - Hierarchical softmax
 - Negative sampling
 - Subsampling frequent words
- 3 Chạy thử
- 4 Tài nguyên tham khảo
- 5 Hỏi đáp và kết thúc

Word2Vec model (1)

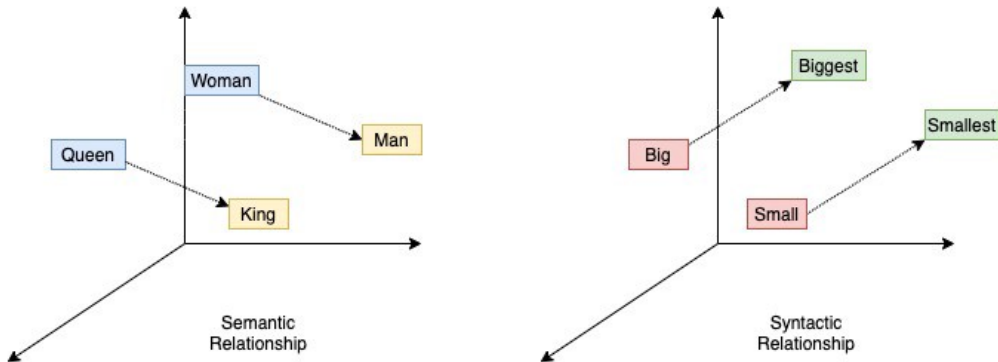


Hình 1: Các từ có ý nghĩa tương đồng nhau nằm gần nhau.

- Word2Vec là một phương pháp biến đổi dữ liệu dạng **text** sang **numeric**.
- Các từ được đại diện bằng các **vector**.
- Word2Vec model sử dụng **neural network** nên yêu cầu input data phải là numeric.
- Khắc phục được nhược điểm của các phương pháp truyền thống như TF-IDF, Bag of Words vì chúng không hiểu được **syntactic** và **semantic** của từ.

Word2Vec model (2)

- Word2Vec là một phương pháp **word embedding** được sử dụng phổ biến.
- Các **word vector** được tạo ra bởi Word2Vec model có khả năng nắm bắt được **semantic** và **syntactic** của từ.



Hình 2: Ví dụ về **semantic** và **syntactic**.

Word2Vec model (3)

- Word2Vec biểu diễn các word vector trong không gian m chiều.
- Ví dụ có câu: *"Archie used to live in New York, he then moved to Santa Clara. He loves apples and strawberries."*



Hình 3: Các từ trong câu ví dụ được biểu diễn trong không gian vector.

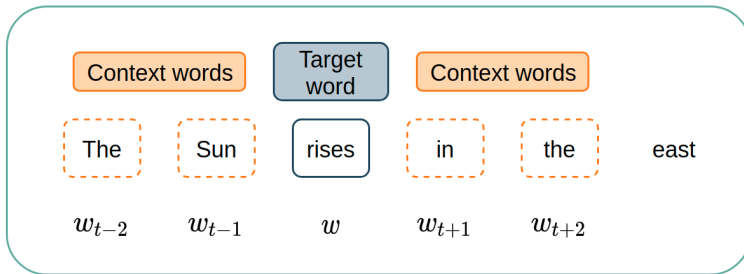
Word2Vec model (4)

- Vì hiểu được semantic và syntactic của từ điều này giúp tận dụng các vector vào các bài toán như **text summarization** (*tóm tắt văn bản*), **sentiment analysis** (*phân tích tình cảm*), **text generation** (*tạo văn bản*),...
- Có hai cách để xây dựng Word2Vec model:
 - ① **CBOW model.**
 - ② **Skip-gram model.**

Word2Vec model (5)

CBOW model

- Giả sử có một neural network bao gồm: **một input layer**, **một hidden layer** và **một output layer**. Cần dự đoán ra **một từ** dựa vào **các từ xung quanh nó**. Từ được dự đoán được gọi là **target word** và các từ xung quanh nó được gọi là **context word**.
- Đặt n là số context word cần thiết để dự đoán target word. Xét $n = 2$ với câu: *"The Sun rises in the east."*

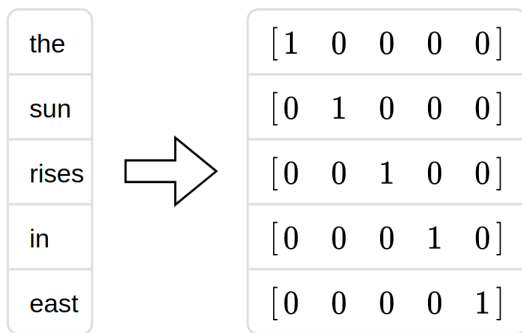


Hình 4: Target word và context word trong CBOW model.

Word2Vec model (6)

CBOW model

- Input của network là các context word và output của network là target word.
- Cần sử dụng kỹ thuật **one-hot encoding** để chuyển đổi các text data thành numeric data trước khi đưa vào network.

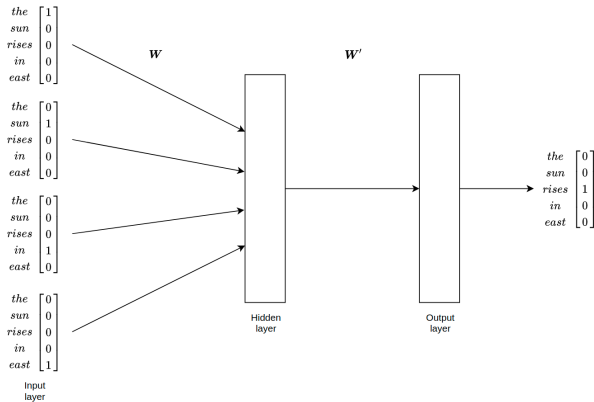


Hình 5: One-hot encoding cho text data.

Word2Vec model (7)

CBOW model

- Kiến trúc của CBOW model. Các context word là: *the*, *sun*, *in* và *east* được dùng làm đầu vào cho network và target word là *rises* được dự đoán ở đầu ra.



Hình 6: Kiến trúc của CBOW model 1.

Word2Vec model (8)

CBOW model

- Kết thúc quá trình đào tạo, lấy weight \mathbf{W} ra làm các word vector cho các vocabulary.
- Dưới đây là các vector tương ứng cho các từ của \mathbf{W} . Word embedding tương ứng cho từ *sun* là $[0.0 \ 0.3 \ 0.3 \ 0.6 \ 0.1]$.

$$\mathbf{W} = \begin{matrix} & \begin{matrix} the \\ sun \\ rises \\ in \\ east \end{matrix} \end{matrix} \begin{bmatrix} 0.01 & 0.02 & 0.1 & 0.5 & 0.37 \\ 0.0 & 0.3 & 0.3 & 0.6 & 0.1 \\ 0.4 & 0.34 & 0.11 & 0.61 & 0.43 \\ 0.1 & 0.11 & 0.1 & 0.17 & 0.369 \\ 0.33 & 0.4 & 0.3 & 0.17 & 0.1 \end{bmatrix}$$

Word2Vec model (9)

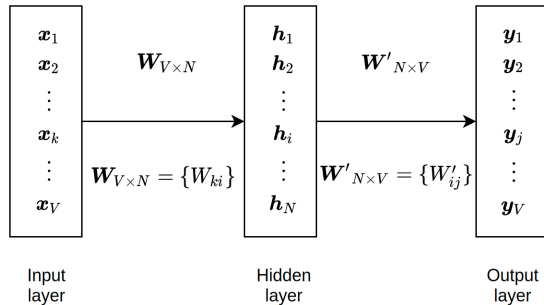
CBOW model - CBOW model với một context word

- Xét trường hợp chỉ sử dụng duy nhất một context word, tức $C = 1$. Lúc này network nhận vào một context word ở đầu vào và trả về một target word ở đầu ra.
- Xét câu: *"The Sun rises in the east."*
- Đặt V là số lượng vocabulary và N là số neuron của hidden layer. Vì neural network của chúng ta có ba layer, cụ thể:
 - Input layer được đại diện bởi $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \dots, \mathbf{x}_k, \dots, \mathbf{x}_V\}$. Khi chúng ta nói \mathbf{x}_k , tức ta đang đề cập đến từ thứ k trong corpus của chúng ta.
 - Hidden layer được đại diện bởi $\mathbf{H} = \{\mathbf{h}_1, \mathbf{h}_2, \mathbf{h}_3, \dots, \mathbf{h}_i, \dots, \mathbf{h}_N\}$. Khi chúng ta nói \mathbf{h}_i , tức ta đang đề cập đến neuron thứ i trong hidden layer.
 - Output layer được đại diện bởi $\mathbf{Y} = \{\mathbf{y}_1, \mathbf{y}_2, \mathbf{y}_3, \dots, \mathbf{y}_j, \dots, \mathbf{y}_V\}$. Khi chúng ta nói \mathbf{y}_j , tức ta đang đề cập đến từ thứ j trong output layer.

Word2Vec model (10)

CBOW model - CBOW model với một context word

- Số chiều của \mathbf{W} là $V \times N$, tức số *vocabulary* \times số *neuron trong hidden layer*. W_{ki} đại diện cho một phần tử trong ma trận \mathbf{W} giữa \mathbf{x}_k trong input layer và \mathbf{h}_i trong hidden layer.
- Số chiều của \mathbf{W}' là $N \times V$, tức số *neuron trong hidden layer* \times số *vocabulary*. W'_{ij} đại diện cho một phần tử trong ma trận \mathbf{W}' giữa \mathbf{h}_i trong hidden layer và \mathbf{y}_j trong output layer.



Hình 7: Kiến trúc của CBOW model 2.

Word2Vec model (11)

CBOW model - CBOW model với một context word

Forward propagation

- **Bước 1:**

$$H = XW^T = W_{(k,.)}$$

$W_{(k,.)}$ là biểu diễn vector của input word. Đặt vector đại diện cho input word w_I bằng Z_{w_I} . Lúc này phương trình trên có thể được viết thành.

$$H = Z_{w_I} \tag{1}$$

Word2Vec model (12)

CBOW model - CBOW model với một context word

Forward propagation

- **Bước 2:** Đặt u_j là điểm số để từ thứ j trong corpus là target word - được tính bằng cách nhân H cho W' .

$$u_j = W'_{ij}^T \cdot H$$

- **Bước 3:** Cột j của W'_{ij} đại diện cho vector của từ j trong corpus. Đặt vector đại diện cho từ thứ j là $Z'_{w_j^T}$

$$u_j = Z'_{w_j^T} \cdot H \quad (2)$$

Word2Vec model (13)

CBOW model - CBOW model với một context word

Forward propagation

- **Bước 4:** Thế phương trình (1) vào phương trình (2), ta được:

$$\mathbf{u}_j = \mathbf{Z}_{w_j'}^T \cdot \mathbf{Z}_{w_I}$$

- **Bước 5:** Áp dụng softmax activation để chuyển đổi \mathbf{u}_j sang xác suất:

$$y_j = \frac{\exp(\mathbf{u}_j)}{\sum_{j'=1}^V \exp(\mathbf{u}_{j'})} \quad (3)$$

Word2Vec model (14)

CBOW model - CBOW model với một context word

Forward propagation

- **Bước 6:** Đặt y_j^* là xác suất để từ j chính xác là target word. Và vì chúng ta áp dụng gradient descent để tìm ra optimal weight nên thay vì tối đa hóa ta sẽ tìm tối thiểu hóa của $\log(y_j^*)$

$$\min(-\log(y_j^*))$$

- **Bước 7:** Loss function của chúng ta lúc này thành:

$$\begin{aligned}\mathcal{L} &= -\log(y_j^*) \\ &= -\log\left(\frac{\exp(u_j)}{\sum_{j'=1}^V \exp(u_{j'})}\right)\end{aligned}$$

Word2Vec model (15)

CBOW model - CBOW model với một context word

Forward propagation (*tiếp theo*)

$$\begin{aligned}\mathcal{L} &= - \left(\log(\exp(u_j)) - \log \left(\sum_{j'=1}^V \exp(u'_{j'}) \right) \right) \\ &= - \log(\exp(u_j)) + \log \left(\sum_{j'=1}^V \exp(u'_{j'}) \right) \\ &= -u_j + \log \left(\sum_{j'=1}^V \exp(u'_{j'}) \right)\end{aligned}$$

Word2Vec model (16)

CBOW model - CBOW model với một context word

Backward propagation

- **Bước 1:** Trong quá trình này, chúng ta sẽ cập nhật cho W và W' như sau:

$$W = W - \alpha \frac{\partial \mathcal{L}}{\partial W}$$

$$W' = W' - \alpha \frac{\partial \mathcal{L}}{\partial W'}$$

Word2Vec model (17)

CBOW model - CBOW model với một context word

Backward propagation

- **Bước 2:** Nhắc lại các công thức trong quá trình forward propagation:

$$H = XW^T$$

$$u_j = W'_{ij}^T \cdot H$$

$$\mathcal{L} = -u_j + \log \left(\sum_{j'=1}^V \exp(u'_{j'}) \right)$$

Word2Vec model (18)

CBOW model - CBOW model với một context word

Backward propagation

- **Bước 3:** Tính gradient của loss function đối với \mathbf{W}' theo quy tắc **chain rule**.

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}'_{ij}} = \frac{\partial \mathcal{L}}{\partial \mathbf{u}_j} \cdot \frac{\partial \mathbf{u}_j}{\partial \mathbf{W}'_{ij}}$$

Ta có đạo hàm của $\frac{\partial \mathcal{L}}{\partial \mathbf{u}_j}$ như sau: $\frac{\partial \mathcal{L}}{\partial \mathbf{u}_j} = \mathbf{e}_j$ (5)

với \mathbf{e}_j là sự khác biệt giữa actual word và predicted word.

Word2Vec model (19)

CBOW model - CBOW model với một context word

Backward propagation

- **Bước 4:** Tính đạo hàm cho $\frac{\partial u_j}{\partial W'_{ij}}$, và bởi vì chúng ta biết

$$u_j = W'_{ij}^T \cdot H \text{ nên:}$$

$$\frac{\partial u_j}{\partial W'_{ij}} = H$$

- **Bước 5:** Như vậy, gradient của loss function đối với W' bằng:

$$\frac{\partial \mathcal{L}}{\partial W'_{ij}} = e_j \cdot H$$

Backward propagation

- **Bước 6:** Tính gradient của loss function đối với \mathbf{W} bằng quy tắc chain rule:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{ki}} = \frac{\partial \mathcal{L}}{\partial \mathbf{h}_i} \cdot \frac{\partial \mathbf{h}_i}{\partial \mathbf{W}_{ki}}$$

- **Bước 7:** Tính đạo hàm của $\frac{\partial \mathcal{L}}{\partial \mathbf{h}_i}$, chúng ta áp dụng chain rule:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{h}_i} = \sum_{j=1}^V \frac{\partial \mathcal{L}}{\partial \mathbf{u}_j} \cdot \frac{\partial \mathbf{u}_j}{\partial \mathbf{h}_i}$$

Word2Vec model (21)

CBOW model - CBOW model với một context word

Backward propagation

- **Bước 8:** Thế phương trình (5) vào phương trình trên, ta được:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{h}_i} = \sum_{j=1}^V \mathbf{e}_j \cdot \frac{\partial \mathbf{u}_j}{\partial \mathbf{h}_i}$$

- **Bước 9:** Vì $\mathbf{u}_j = \mathbf{W}'_{ij}^T \cdot \mathbf{H}$, tiếp tục thế vào phương trình trên được:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{h}_i} = \sum_{j=1}^V \mathbf{e}_j \cdot \mathbf{W}'_{ij} = \mathcal{L} \mathbf{H}^T$$

Word2Vec model (22)

CBOW model - CBOW model với một context word

Backward propagation

- **Bước 10:** Tính đạo hàm cho $\frac{\partial h_i}{\partial W_{ki}}$. Biết $H = XW^T$ nên:

$$\frac{\partial h_i}{\partial W_{ki}} = X$$

- **Bước 11:** Cuối cùng, gradient của loss function đối với W là:

$$\frac{\partial \mathcal{L}}{\partial W_{ki}} = LH^T \cdot X$$

Word2Vec model (23)

CBOW model - CBOW model với một context word

Backward propagation

- **Bước 12:** Như vậy, các set weight W và W' được cập nhật trong quá trình backward propagation như sau:

$$W = W - \alpha L H^T \cdot X$$

$$W' = W' - \alpha e_j \cdot H$$

Word2Vec model (24)

CBOW model - CBOW model với một context word

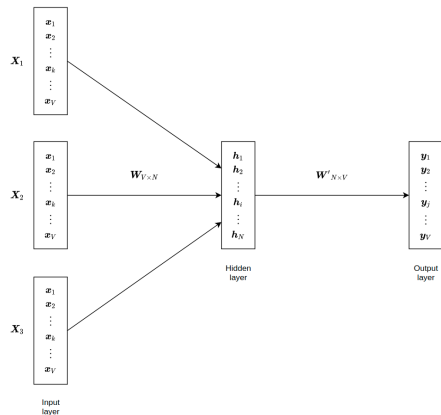
Mã nguồn

```
1 def single_context_CBOW(x, label, W1, W2, loss):
2     # forward propagation
3     h = np.dot(W1.T, x)
4     u = np.dot(W2.T, h)
5     y_pred = softmax(u)
6
7     # error giữa actual value và predicted value
8     dW2 = np.outer(h, e)
9     dW1 = np.outer(x, np.dot(W2.T, e))
10
11    # cập nhật các weight set
12    W1 = W1 - lr * dW1
13    W2 = W2 - lr * dW2
14
15    # loss function
16    loss += -float(u[label] == 1) + np.log(np.sum(np.exp(u)))
17
18    return W1, W2, loss
```

Word2Vec model (25)

CBOW model - CBOW model với multiple context words

- Dưới đây là kiến trúc của CBOW model với multiple context words:



Hình 8: Kiến trúc của CBOW model với multiple context words.

Word2Vec model (26)

CBOW model - CBOW model với multiple context words

- Với multiple context words, chúng ta lấy giá trị trung bình của tất cả các context word đầu vào. Cụ thể với C context word $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_C$, chúng ta thực hiện:

$$\begin{aligned} \mathbf{H} &= \frac{(\mathbf{x}_1 + \mathbf{x}_2 + \dots + \mathbf{x}_C)}{C} \mathbf{w}^T \\ &= \frac{1}{C} (\mathbf{x}_1 \mathbf{w}^T + \mathbf{x}_2 \mathbf{w}^T + \dots + \mathbf{x}_C \mathbf{w}^T) \end{aligned}$$

- Tương tự như CBOW model với single context word, $\mathbf{x}_1 \mathbf{w}^T$ đại diện cho vector của input context word w_1 . Đại diện Z_{w_1} cho input context word w_1 , cho nên:

$$\mathbf{H} = \frac{1}{C} (Z_{w_1} + Z_{w_2} + \dots + Z_{w_C}) \quad (6)$$

Word2Vec model (27)

CBOW model - CBOW model với multiple context words

- Ở đây, C đại diện cho số lượng context word, để tính toán giá trị của \mathbf{u}_j sẽ tương tự như cách ta tính toán cho single context word.

$$\mathbf{u}_j = Z'_{w_j^T} \cdot \mathbf{H} \quad (7)$$

với $Z'_{w_j^T}$ là vector đại diện cho từ thứ j trong corpus.

- Thế phương trình (6) vào phương trình (7), ta được:

$$\mathbf{u}_j = Z'_{w_j^T} \cdot \frac{1}{C} (Z_{w_1} + Z_{w_2} + \dots + Z_{w_C})$$

- Tính toán loss function:

$$\mathcal{L} = -u_j + \log \left(\sum_{j'=1}^V \exp(u_{j'}) \right)$$

Word2Vec model (28)

CBOW model - CBOW model với multiple context words

- Tiếp theo, là quá trình backward propagation. Cụ thể, quá trình cập nhật cho \mathbf{W} sẽ như sau:

$$\mathbf{W} = \mathbf{W} - \alpha \mathbf{L} \mathbf{H}^T \cdot \frac{(\mathbf{x}_1 + \mathbf{x}_2 + \dots + \mathbf{x}_C)}{C}$$

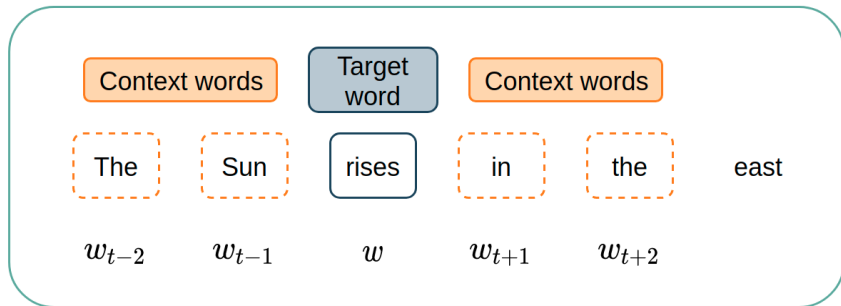
- Đối với \mathbf{W}' sẽ tương tự như single context word:

$$\mathbf{W}' = \mathbf{W}' - \alpha \mathbf{e}_j \cdot \mathbf{H}$$

Word2Vec model (29)

Skip-gram model

- Skip-gram model về cơ bản như là việc làm ngược lại CBOW model. Với Skip-gram model, chúng ta sẽ cố gắng dự đoán các context word dựa vào target word được cung cấp ở đầu vào. Lấy lại ví dụ ban đầu của CBOW model, chúng ta sẽ dùng target word là *rises* để dự đoán các context word là *the*, *sun*, *in* và *east*.

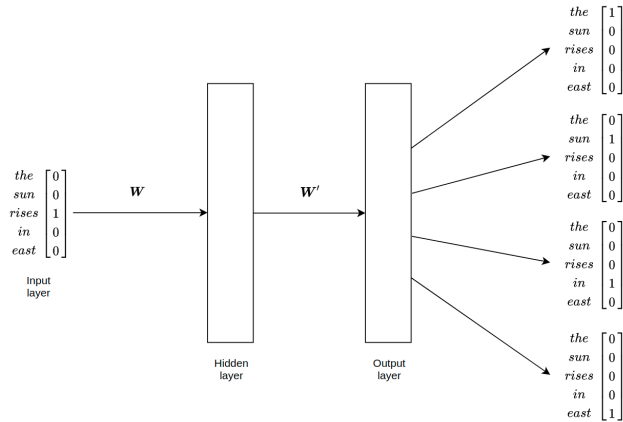


Hình 9: Target word và context word trong Skip-gram model.

Word2Vec model (30)

Skip-gram model

- Dưới đây là kiến trúc của Skip-gram model:

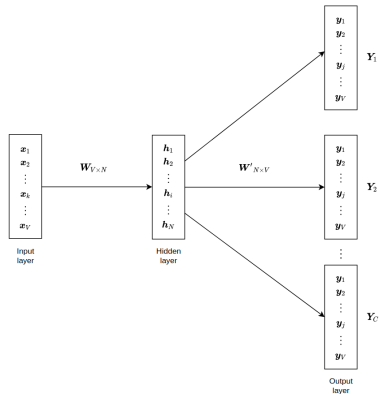


Hình 10: Kiến trúc của Skip-gram model 1.

Word2Vec model (31)

Skip-gram model

- Với một target word làm đầu vào \mathbf{X} và trả về C context word \mathbf{Y} ở đầu ra.



Hình 11: Kiến trúc của Skip-gram model 2.

Forward propagation

- **Bước 1:** Trước tiên chúng ta nhân X cho hidden layer weight W .

$$H = XW^T = Z_{w_l}$$

ở đây, Z_{w_l} là vector đại diện cho input word w_l .

- **Bước 2:** Tính u_j , đây chính là điểm số dùng để đánh giá sự tương đồng giữa từ j trong corpus và target word.

$$\begin{aligned} u_j &= W'_{ij}^T \cdot H \\ &= Z'_{w_j}^T \cdot H \end{aligned}$$

với $Z'_{w_j}^T$ là vector đại diện cho từ j .

Forward propagation

- **Bước 3:** Không giống với CBOW model chỉ cần dự đoán duy nhất một target word. Ở đây, dự đoán ra C context word. Nên để cho tường minh chúng ta có thể ghi lại phương trình trên thành:

$$\mathbf{u}_{c,j} = \mathbf{Z}'_{w_j T} \cdot \mathbf{H} \text{ với } c = 1, 2, 3, \dots, C$$

Vì $\mathbf{u}_{c,j}$ là điểm số đại diện cho từ thứ j là context word c . Cho nên:

- $\mathbf{u}_{1,j}$ là điểm số cho từ thứ j sẽ là context word đầu tiên.
- $\mathbf{u}_{2,j}$ là điểm số cho từ thứ j sẽ là context word thứ hai.
- $\mathbf{u}_{c,j}$ là điểm số cho từ thứ j sẽ là context word thứ c .

Forward propagation

- **Bước 4:** Chuyển toàn bộ các điểm số trên thành xác suất bằng hàm softmax và đại diện bởi $\mathbf{y}_{c,j}$.

$$\mathbf{y}_{c,j} = \frac{\exp(\mathbf{u}_{c,j})}{\sum_{j'=1}^V \exp(\mathbf{u}_{j'})} \quad (8)$$

$\mathbf{y}_{c,j}$ là xác suất cho từ thứ j trong corpus sẽ là context word thứ c .

Forward propagation

- **Bước 5:** Đặt $y_{c,j}^*$ là xác suất để từ được chọn đúng là context word. Chúng ta cần tối thiểu hóa hàm $-\log$ xác suất này.

$$\min(-\log(y_{c,j}^*))$$

- **Bước 6:** Thế phương trình (8) vào phương trình trên ta được:

$$\mathcal{L} = -\log \frac{\exp(u_{c,j})}{\sum_{j'=1}^V \exp(u_{j'})}$$

Forward propagation - (tiếp theo)

Và vì chúng ta có C context word, nên chúng ta sẽ lấy tích các xác suất này với nhau.

$$\begin{aligned}\mathcal{L} &= -\log \prod_{c=1}^C \frac{\exp(\mathbf{u}_{c,j})}{\sum_{j'=1}^V \exp(\mathbf{u}_{j'})} \\ &= -\sum_{c=1}^C \mathbf{u}_{c,j} + C \cdot \log \sum_{j'=1}^V \exp(\mathbf{u}_{j'})\end{aligned}$$

Backward propagation

- **Bước 1:** Áp dụng gradient descent, tính toán gradient cho loss function và cập nhật cho weight. Trước tiên, tính toán gradient cho W' bằng chain rule như sau:

$$\frac{\partial \mathcal{L}}{\partial W'_{ij}} = \sum_{c=1}^C \frac{\partial \mathcal{L}}{\partial \mathbf{u}_{c,j}} \cdot \frac{\partial \mathbf{u}_{c,j}}{\partial W'_{ij}}$$

- **Bước 2:** Đạo hàm của $\frac{\partial \mathcal{L}}{\partial \mathbf{u}_{c,j}}$ là sự khác biệt giữa actual value và predicted value: $\frac{\partial \mathcal{L}}{\partial \mathbf{u}_j} = \mathbf{e}_{c,j}$

Backward propagation

- **Bước 3:** Do có nhiều context word, nên cộng tất cả các $\mathbf{e}_{c,j}$ này lại:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{u}_j} = \sum_{c=1}^C \mathbf{e}_{c,j} = \mathbf{E}$$

- **Bước 4:** Tính đạo hàm cho $\frac{\partial \mathbf{u}_j}{\partial \mathbf{W}'_{ij}}$, và bởi vì chúng ta biết

$$\mathbf{u}_j = \mathbf{W}'_{ij}^T \cdot \mathbf{H}, \text{ cho nên:}$$

$$\frac{\partial \mathbf{u}_j}{\partial \mathbf{W}'_{ij}} = \mathbf{H}$$

Backward propagation

- **Bước 5:** Cho nên, gradient của loss function đối với \mathbf{W}' như sau:

$$\frac{\partial u_j}{\partial \mathbf{W}'_{ij}} = E \cdot H$$

- **Bước 6:** Tính gradient của loss function đối với \mathbf{W} .

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{ki}} = \frac{\partial \mathcal{L}}{\partial \mathbf{h}_i} \cdot \frac{\partial \mathbf{h}_i}{\partial \mathbf{W}_{ki}}$$

Backward propagation

- **Bước 7:** Lúc này, gradient của loss function so với W và W' như sau:

$$W' = W' - \alpha E \cdot H$$

$$W = W - \alpha L H^T \cdot X$$

- **Bước 8:** Kết thúc quá trình đào tạo, chúng ta sẽ tiến hành cập nhật cho các weight set trên network. Sau cùng W đạt trạng thái optimal weight và trở thành word vector cho vocabulary trong corpus.

Các chiến lược cải thiện hiệu suất (1)

Để tiến hành cải thiện hiệu suất của quá trình đào tạo model, người ta sẽ sử dụng kèm theo một số kỹ thuật như:

- Sử dụng **Hierarchical softmax** thay vì softmax truyền thống.
- **Negative sampling**.
- **Subsampling frequent words**.

Các chiến lược cải thiện hiệu suất (2)

Hierarchical softmax

- Tính toán bằng softmax sẽ rất tốn kém về mặt hiệu suất nên cần sử dụng một cách khác để tối ưu hóa hiệu năng tính toán đó là tính softmax trên cây - Hierarchical softmax:

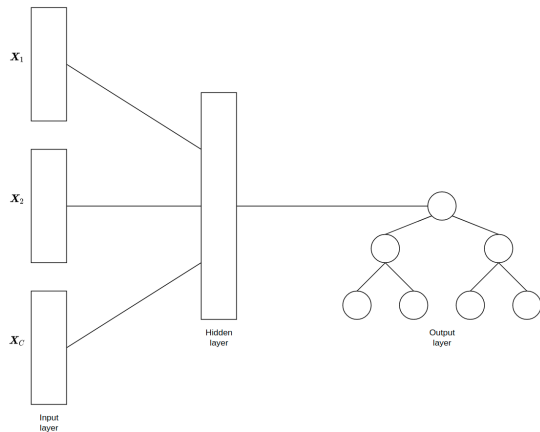
$$y_j = \frac{\exp(u_j)}{\sum_{j'=1}^V \exp(u_{j'})}$$

độ phức tạp $O(V)$.

- Hierarchical softmax sử dụng **Huffman binary search tree** (cây tìm kiếm nhị phân *Huffman*) để giảm độ phức tạp về mặt tính toán xuống $O(\log_2(V))$. Lúc này kiến trúc mạng có sự thay đổi đáng kể ở output layer - chúng ta sẽ thay output layer bằng một **binary search tree**.

Các chiến lược cải thiện hiệu suất (3)

Hierarchical softmax

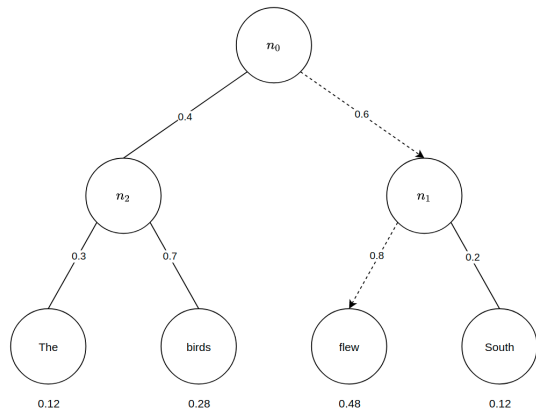


- Đây là kiến trúc của Hierarchical softmax sau khi được nhúng vào output layer của network.
- Mỗi node lá bây giờ trong cây đại diện cho một từ trong corpus và tất cả các node trung gian đại diện cho **relative probability** xác suất tương đối của tất cả các node con của chúng.

Hình 12: Thay thế output layer bằng Hierarchical softmax trong CBOW model.

Các chiến lược cải thiện hiệu suất (4)

Hierarchical softmax



Hình 13: Thay thế output layer bằng Hierarchical softmax trong CBOW model.

- Để tính xác suất của một target word khi được cung cấp context word. Chỉ cần duyệt cây và đưa ra quyết định nên duyệt cây bên trái hay bên phải tại một node nhất định. Giả sử cần tính xác suất của từ *flew* là target word dựa vào *c* context word. Chúng ta lấy tích của tất cả các xác suất dọc theo con đường đến target word.

$$\begin{aligned} p(flew|c) &= p_{n_0}(left|c) \times p_{n_1}(right|c) \\ &= 0.6 \times 0.8 \\ &= 0.12 \end{aligned}$$

Các chiến lược cải thiện hiệu suất (5)

Negative sampling

- Xét câu: "*Birds are flying in the sky.*" với target word là *flying* và còn lại là các context word. Chúng ta cần cập nhật lại các weight trong network mỗi khi nó dự đoán một target word không chính xác. Vì vậy, ngoại trừ từ *flying*, nếu một từ khác được dự đoán là target word thì chúng ta sẽ cập nhật network.
- Như vậy, với corpus chứa hàng triệu vocabulary sẽ rất hao tổn tài nguyên về mặt tính toán.
- Để khắc phục điều này, đánh dấu correct target word là positive class và lấy mẫu một vài từ còn lại thuộc negative class. Về cơ bản, chúng ta đang chuyển về bài toán phân loại nhị phân. Lúc này, xác suất để một từ được chọn là negative word.

$$p(w_i) = \frac{\text{frequency}(w_i)^{\frac{3}{4}}}{\sum_{j=0}^n \text{frequency}(w_j)^{\frac{3}{4}}}$$

Các chiến lược cải thiện hiệu suất (6)

Subsampling frequent words

- Trong corpus, sẽ có một số từ xuất hiện rất thường xuyên chẳng hạn như *the*, *is*,... và có một số từ lại xuất hiện không thường xuyên cho lắm. Để duy trì sự cân bằng giữa hai yếu tố, người ta sử dụng kỹ thuật **subsampling** (lấy mẫu con). Vậy nên, chúng ta sẽ loại bỏ một số từ xuất hiện thường xuyên hơn một **threshold** nhất định với xác suất p được tính như sau:

$$p(w_i) = 1 - \sqrt{\frac{t}{f(w_i)}}$$

ở đây t là threshold và $f(w_i)$ là frequency của từ i .



- **Hands-On Deep Learning Algorithms with Python** - (chương 7), link sách tại <https://www.amazon.com/Hands-Deep-Learning-Algorithms-Python/dp/1789344158>.
- **Natural Language Processing in Action**, link sách tại <https://www.amazon.com/Natural-Language-Processing-Action-Understanding/dp/1617294632>.



HẾT

Cảm ơn Thầy và các bạn đã chú ý lắng nghe.