

ĐẠI HỌC QUỐC GIA TP. HỒ CHÍ MINH
TRƯỜNG ĐẠI HỌC CÔNG NGHỆ THÔNG TIN



BÀI TẬP

MÔN CÁC KỸ THUẬT HỌC SÂU

VÀ ỨNG DỤNG

KHOA: KHOA HỌC MÁY TÍNH

HOMEWORK: GIẢI THÍCH HÀM HIERALCHICAL SOFTMAX

Giảng viên: Nguyễn Duy Khánh

Nhóm thực hiện:

1. Trương Văn Khải– 21520274

Hình 1

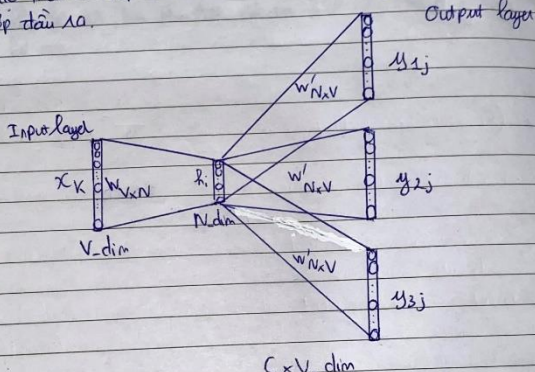
Paper

Date

Giải thích Hierarchical Softmax

* Mô hình Skip-gram

- Mô hình này hoạt động ngược với mô hình CBOW. Từ ngữ mục tiêu ở lớp đầu vào và các từ cùng ngữ cảnh sẽ ở lớp đầu ra.



- Mục tiêu của skip-gram: tìm từ đại diện để dự đoán các từ xung quanh trong 1 câu hay 1 từ hiện. Ngoài ra, nó còn có nhiệm vụ huấn luyện 1 tập $w_1, w_2, \dots, w_t \Rightarrow$ Tập
- \Rightarrow Mục tiêu của skip-gram là tối ưu hóa 2 xác suất trung bình log, ta có công thức:

$$\frac{1}{T} \sum_{t=1}^T \log p(w_{t+j} | w_t)$$

$$t=1-c \leq j \leq c, j \neq 0$$

Trong đó, c là quy mô của ngữ cảnh huấn luyện (thông thường có thể là một hàm số tăng dần w_t), và việc xây dựng skip-gram là để bài toán tính $p(w_{t+j} | w_t)$ bằng cách sử dụng hàm softmax:

QT BOOK

Hình 2

Output Layer

đơn các
giáo ra, nó
t → T₀
suất tung

nguyên (đồng
và việc
t + j | w_t)

Paper:
Date:

$$p(w_t | w_{t-1}) = \frac{\exp(v' w_t)}{\sum_{w=1}^W \exp(v' w_t)}$$

Trong công thức trên, v_w và $v_{w'}$ là các vector "đầu vào" và "đầu ra" của w và w' là số từ tương tự điển.

Nhưng công thức này không thực tế vì giá trị của phép tính $\sum \log p(w_t | w_{t-1})$ tỷ lệ thuận với W , mà giá trị này rất lớn (từ 10⁵ đến 10⁷).

*** Hierarchical Softmax (Softmax phân cấp)**

- Vì vậy, ta cần tìm một phép tính xấp xỉ nhưng hiệu quả ngang ngửa hàm Softmax. Hierarchical Softmax có thể giải quyết nhiệm vụ này.
- Đầu vào mỗi nút, Hierarchical Softmax sử dụng 1 cây đại diện nhị phân của lớp đầu ra với các từ w như là lá của nó. Đồng thời, 1 đường đi ngẫu nhiên được cho là xác suất đầu vào các từ.

Softmax Phân Cấp:
Mỗi nút lá của cây biểu diễn 1 từ trong từ điển.

Rõ ràng hơn, mỗi từ w có thể đạt được bằng một đường từ gốc của cây. Gọi $n(w, j)$ là node thứ j

QTBOOK

Hình 3

