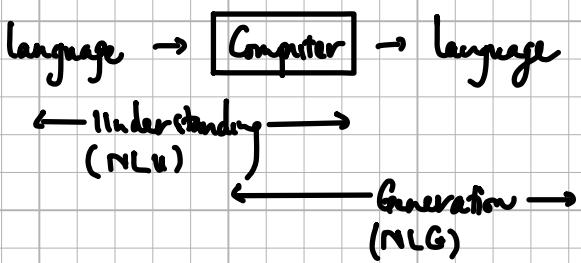


Viet AI - NLP

Week 1 : NLP and word embeddings

- Goals:



- Deep Learning for NLP:

- Recurrent neural networks
- Long short-term memory
- Convolutional neural networks
- Autoencoders
- Transformer

- Vector Representation of Words

word → \vec{w} / vector → +, -, *, /

↳ One-hot vector, Bag-of-words or BOW (1)

↳ Contextual representation (2)

(1) BOW: \vec{w}_t : $\text{dim} \in \mathbb{R}^n$ Số lượng từ \rightarrow sparse matrix (Chữ và trống)
 Thuyết: tính độ tần số

(2) Contextual Representation

- word vector
- using dot product

Word2Vec :

↳ Skip gram: $L(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{-m \leq j \leq m} \log P(w_{t+j} | w_t; \theta)$

↳ CBOW: $L(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{j \neq 0} \log P(w_t | w_{t-m}, \dots, w_{t+m}; \theta)$

- 2.7 + D2L Prob
- Machine B2Sán
+ Kèm file VietAI
- Ielts Reading
+ Listening
- Ielts Writing
- Coursera
- Personal Projects

RNN Variants

Đoạn ký tự chuỗi, từ ảnh, feature

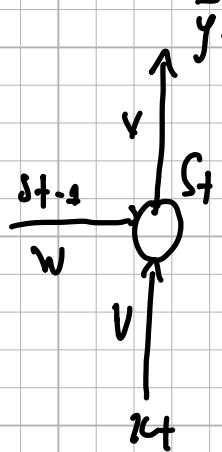
- Tính thứ tự

Sequence data Processing

↳ Recurrent Network

unfolded form

→ input [ký] → hidden [tự] → output



thông
tự

s_t

$$s_t = \sigma(Vx_t + Ws_{t-1})$$

wtai

SCR

$$\tilde{y}_t = \text{softmax}(Vs_t)$$

distribution

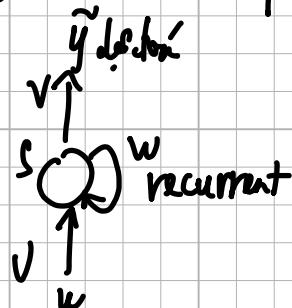
decide/选出

hà

Tóm tắt RNN không bao hàm

U, V, W

chưa ghi các time step



→ Compressed form

[ký] → phi tuyến

[Step 1]

Model Design

$$\text{sigmoid: } \frac{1}{1+e^{-u}}$$

$$\text{softmax: } \frac{e^u}{\sum e^u}$$

[Step 2] loss function Design

$$L(\theta) = - \sum_{j=1}^k y_{t,j} \log \tilde{y}_{t,j}$$

with $y_{t,j}$

$$L(\theta) = - \frac{1}{T} \sum_{t=1}^T L_t(\theta)$$

Learn this: (func' LSTM, GRU

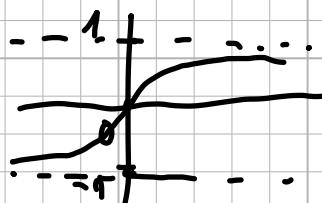
Some Use cases:

- ↳ one to one , one to many , many to one
many to many (1) , many to many (2)

• Training process : loss function

⇒ tanh function

- Cell \Rightarrow sigmoid stalk.
- $H^t = W \cdot H^{t-1}$
- LSTM



Ablation Study: which was loss sig

• Problem of Parallel computing

HCMUS

- [] Prod and Stats
- [] train h_i t² hisp

[] 27 + D2L Prob + D2L 246

[] Machine Translation

+ Bao tap ViennAI

[] IELTS Reading + Writing
+ Listening + Speaking

Certi

[] CSE229

[] IELTS Writing

[] DL

[] Coursera

[] NLP

[] Personal Projects

[] RL

[] 100 DS

Attention(Q, K, V)

$$= \text{softmax} \left(\frac{QK^T}{\sqrt{dk}} \right) \cdot V$$

Attention is all you need

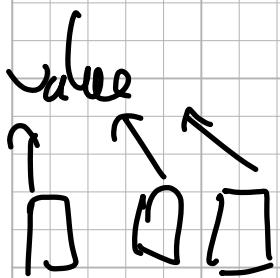
• Hidden states (key , query , value)

• Self Attention

• Transformer Architecture

the self Attention \rightarrow FFNN

Key, query, value



h...
query
key

• Multi-head Attention

↳ FFNN Concat and shorten ?

Ex 8.

- D2L Prob + D2L C3, 1 3 5 7

Estimation of the mean

→ Normal IID samples

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$$

$$E(\bar{X}_n) = E\left[\frac{1}{n} \sum_{i=1}^n X_i\right] = \frac{1}{n} \cdot n\mu = \mu$$

$$\text{Var}[\bar{X}_n] = \frac{1}{n^2} \cdot n\sigma^2 = \frac{\sigma^2}{n}$$

$$\bar{X}_n \sim N(\mu, \frac{\sigma^2}{n})$$

$$P=1$$

Ex 3: Prob we see a head after drawing N exams

$$\textcircled{1} \quad P_N(p) = C_N^n \frac{1}{2^n} \cdot \frac{1}{2^{N-n}} = C_N^n \frac{1}{2^N} \quad (n: 0..N)$$

$$\begin{aligned} E_N[\tilde{p}_N] &= \sum_{n=1}^N \frac{n}{N} C_N^n \cdot \frac{1}{2^n} \\ &= \frac{1}{2^N} \left(\sum_{n=1}^N \frac{n}{N} \frac{N!}{(N-n)!n!} \right) = \frac{1}{2} \end{aligned}$$

$$\sqrt{N}[m] = N \cdot p \cdot (1-p) = \frac{N}{4}$$

$$\text{Var}\left[\frac{m}{N}\right] = \frac{1}{N^2} \cdot \frac{N}{4} = \frac{1}{4N}$$

\textcircled{2} Chebychev Inequality (Deviation $f > 0$)

$$P_N\left(|\tilde{p} - \frac{1}{2}| \geq f\right) \leq \frac{1}{4N} f^2$$

$$\textcircled{3} \quad \bar{X} \sim N(\mu, \frac{\sigma^2}{n}) \quad \tilde{p} \sim N\left(\frac{1}{2}, \frac{1}{4N}\right)$$

Ex 4: $x_i \sim N(0, 1)$

Compute $z_m = \frac{1}{m} \sum_{i=1}^m x_i$

$\cdot E(z_m) = E\left(\frac{1}{m} \cdot \sum_{i=1}^m x_i\right) = \frac{1}{m} \cdot E\left(\sum_{i=1}^m x_i\right) = \frac{1}{m} \cdot 0 = 0$

$\text{Var}(z_m) = \frac{1}{m^2} \cdot \sum_{i=1}^m \text{Var}(x_i) = \frac{1}{m^2} \cdot m \cdot 1 = \frac{1}{m}$

- Cannot apply Chebyshev's inequality for every z_m independently.
→ Not independent events, Dependence on sample size,
Collective behavior

Ex 5: $P(A), P(B)$

$$0 \leq P(A \cup B) \leq P(A) + P(B)$$

$$0 \leq P(A \cap B) \leq \min(P(A), P(B))$$

Ex 6: Markov Chain, B only depends on A, C only depends on B

$$P(A|BC) = P(A) \cdot P(B|A) \cdot P(C|B)$$

Ex 7: Assume two test are not independent

$$P(D_1=1 | H=0) = 0,1 \text{ and } P(D_2=0 | H=1) = 0,9$$

$$\text{if } H=1 : P(D_1, D_2 | H=1) = P(D_1 | H=1) \cdot P(D_2 | H=1)$$

$$\text{and } P(D_1 = D_2 = 1 | H=0) = 0,02$$

$$\begin{array}{ccccc} H=1 & & H=0 & & \\ P(D_1=1 | H) & 0,99 & 0,1 & P(D_2=1 | H) & 0,99 & 0,1 \\ P(D_1=0 | H) & 0,01 & 0,9 & P(D_2=0 | H) & 0,01 & 0,9 \end{array}$$

$$P(D_1, D_2 | H=0) \cdot P(H=0) (= P(D_1, D_2, H=0))$$

$$= P(D_2 | D_1, H=0) \cdot P(D_1 | H=0) \cdot P(H=0)$$

	$H=1$	$H=0$	$P(D_2=1 D_1=0, H=0)$ $= 0,08$
$P(D_1=D_2=1 H)$	0,5801	0,02	
$P(D_1=1, D_2=0 H)$	0,0095	0,08	YES!
$P(D_1=0, D_2=1 H)$	0,0095	0,08	$P(D_2=0 D_1=0, H=0)$
$P(D_1=D_2=0 H)$	0,50001	0,82	$\rightarrow 0,82 0,02 = 0,91$

$$(2) P(H=1) = 0,0015$$

$$P(H=1 | D_1=1) = \frac{P(D_1=1 | H=1) \cdot P(H=1)}{P(D_1=1)}$$

↙

$$P(D_1=1, H=1) + P(D_1=1, H=0)$$

$$= P(D_1=1 | H=1) \cdot P(H=1) + P(D_1=1 | H=0) \cdot P(H=0)$$

$$(3) P(H=1 | D_1=1, D_2=1)$$

$$= \frac{P(D_1=1, D_2=1 | H=1) \cdot P(H=1)}{P(D_1=1, D_2=1)}$$

↙

$$P(D_1=1, D_2=1 | H=1) \cdot P(H=1) + P(D_1=1, D_2=1 | H=0) \cdot P(H=0)$$

$$(*) = P(D_1=1 | D_2=0, H=0) \cdot P(D_2=0 | H=0)$$

$$\text{Tai: } D_1 = D_2 = 1$$

$$\Rightarrow P(D_1=1 | D_2=0, H=0) = 0,08$$

$$0,02 = P(D_2=1 | D_1=1, H=0) \quad P(D_1=0 | D_2=0, H=0) = 0,5 \quad (*)$$

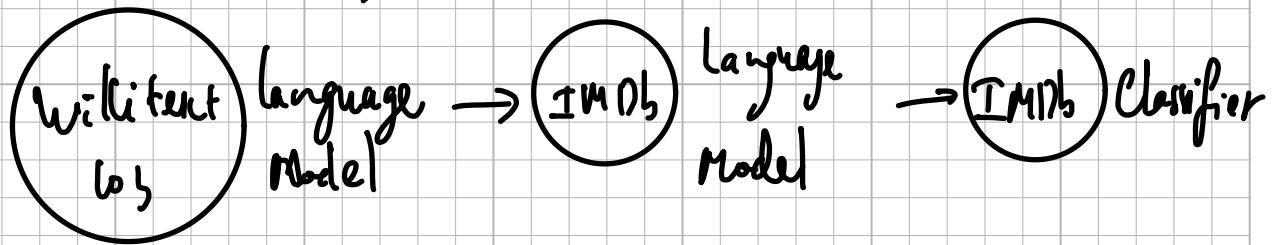
$0,1$

$$\Rightarrow P(D_2=1 | D_1=1, H=0) = 0,2$$

$$\Rightarrow P(D_2=0 | D_1=1, H=0) = 0,8 \quad \Rightarrow P(D_1=1, D_2=0 | H=0) \quad (*)$$

[Coding 2.1] Using pre-trained Datasets and Text Summarization

Transfer learning → Multi-task learning



bert: Bidirectional Encoder

Input Representation

- Workpiece: tokenizer / [CLS]

FFINN + softmax

→ masked language Model

why? Too little: too expensive

Too much: Not enough context

→ Next sentence prediction (Not commonly used)

- Byte-pair Encoding
Unknown token

- Long texts

+) Compact Pre-trained Model

Downstream tasks

↳ sequence Classification

↳ Token classification

+) Token classification

4 last Bert

+) Fine - tuning Strategy

Adapter LORA

→ Word tokenization

+ <UNK>

→ sub word tokenization (units such as morphemes or phonemes)

1994 pair encoding

→ byte-pair encoding

1.1.1. Exercises

1. Data $x_1, \dots, x_n \in \mathbb{R}$. find constant b

$$\underset{b}{\operatorname{Argmin}} \sum_i (x_i - b)^2$$

$$(1) f(b) = \sum_i (x_i - b)^2 \Rightarrow f' = - \sum_i 2(x_i - b) = 0$$

$$\Rightarrow \sum_i x_i = nb \Rightarrow b = \frac{\sum_i x_i}{n}$$

$f'' > 0 \Rightarrow$ convex \Rightarrow

$$\min_b (x - b)^\top (x - b)$$

$$((x_i - b)^2)' = \sum_i 2(x_i - b) = 0$$

$$\sum_i x_i = nb.$$