

1. The course

- technics for NLP
- A big picture understanding of human languages + 구현이 어려운 이유
- PyTorch를 이용한 NLP의 주요 문제 이해 + 구현

2. Human language and word meaning

언어의 사후성 → 컴퓨터에게 어려운 이유

인간 언어는 굉장히 추상적 개념. 부족한 다른 구성원과의 소통을 원활히 해줌

쓰레기 발발 → 거리를 영어서 소름

GPT-3가 하는 일은 다음에 올 단어를 찾는 것. SQL 같은 언어도 유한

언어를 컴퓨터에게 이해시키기

- Common NLP solution: WordNet (유의어와 상위어가 한 set에 있는 사전) 같은 thesaurus (유의어 사전)를 만들
→ 한계: 상황에 따른 차이를 볼 수 없음
단어의 새로운 의미는 // → 추측하기 위해선 도움이 많이 필요

traditional NLP의 문제

- 단어가 별개의 symbol로 여겨짐

one-hot encoding을 사용하면 단어 하나당 하나의 벡터로 바꿔주기 때문에 Vector dimension이 너무 커짐
= 단어의 수

게다가 vector끼리의 similarity를 표현할 방법 X

→ (과거의 시도) WordNet의 유의어 리스트를 이용해보려 했으나 실패

→ vector 자체에 similarity를 담은 방법을 배워보자

< Distributional Semantics >

'단어의 의미는 자주 등장 내에서 가까이 쓰이는 단어들로 결정된다' 는 아이디어에서 출발

banking = $\begin{bmatrix} 1.2 \\ 0.3 \\ 2.5 \end{bmatrix}$ 처럼 n-dimensional vector로 나타냄
→ 실제로는 약 300

이런 Word Vectors는 Word Embeddings. (neural) word representations 라고도 부름

유의어들끼리 가까운 거리에 있도록 설정

이게 바로 word2vec!

3. Word2vec Introduction

center word c가 주어졌을 때 context ('outside') words o가 문장 내에 나올 확률을 알아냄

계산된 p를 이용해 c의 vector 값에 맞춰 o의 vector를 정함

그러고 o를 center word로 바꾸고 그 옆 단어를 context word로 정함

4. Word2vec objective function gradients

* objective function = cost J = loss J

Likelihood = $L(\theta) = \prod_{t=1}^T \prod_{j=1}^M p(w_{t+j} | w_t; \theta)$ conditional probability

objective function = $J(\theta) = -\frac{1}{T} \log L(\theta)$

goal: maximize the likelihood (= minimize the objective function)

응답 예측을 위해 제곱 dot product a similarity
$$P(o|c) = \frac{\exp(u_o \cdot v_c)}{\sum_{w \in V} \exp(u_w \cdot v_c)}$$

 $u \cdot v = u \cdot v = \sum_{i=1}^n u_i \cdot v_i$

$\mathbb{R}^n \rightarrow (0,1)^n$ 인 softmax function가 바로 이것! (scaling해준다는 뜻)
제곱해서 큰 수를 강조하게 됨. max의 역할

* max는 최댓값을 반환하지만 softmax는 점진적으로 그 자체를 반환

5. Optimization basics

$\theta = \begin{bmatrix} w_1 \\ \vdots \\ w_n \end{bmatrix} \in \mathbb{R}^{2d \times V}$ 가장자리 모든 모든 parameter가 담긴 하나의 긴 vector

우리가 배우는 모델에서 θ 는 word vector 하나 뿐

* 모든 단어는 center vector와 context vector 2개를 가진다



결과 수학적으로 모든 vector gradients를 계산해 objective J를 최소화해야 함

이유: 하나만 있는 한 단어를 center vector와 context vector로 동시에 쓸 때 미묘이듯 생겨짐

계속 계산하다 보면 두 벡터는 비슷해짐 (동일 X)

< Gradient >

$$J(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \log P(w_{t+j} | w_t)$$

$$p(o|c) = \frac{\exp(u_o^T v_c)}{\sum_{w=1}^V \exp(u_w^T v_c)}$$

$$* \exp(n) = e^n$$

$$\frac{\partial}{\partial v_c} \log \frac{\exp(u_0^T v_c)}{\sum_{w=1}^V \exp(u_w^T v_c)} = \frac{\partial}{\partial v_c} \log \exp(u_0^T v_c) - \frac{\partial}{\partial v_c} \log \sum_{w=1}^V \exp(u_w^T v_c)$$

$$\textcircled{1} \quad \frac{\partial}{\partial \mathbf{V}_c} u_o^T \mathbf{V}_c = \frac{\partial}{\partial \mathbf{V}_c} (u_{o_1} v_{c_1} + u_{o_2} v_{c_2} + \dots + u_{o_d} v_{c_d})$$

\uparrow
 vector

$$= u_{o_1} + u_{o_2} + \dots + u_{o_d} = u_o$$

④ $\frac{\partial}{\partial v_c} \log \sum_{w=1}^K \exp(\kappa_w^T v_c)$ $f(g(v_c))$ 에 chain rule 을 이용할 것 $h(x) = f(g(x))$
 $\frac{f}{f} \frac{z = g(v_c)}{z = g(v_c)}$ $f, z = g(v_c)$ 각각 미분 $h'(x) = f'(g(x))g'(x)$

$$= \frac{\partial}{\partial \mathbf{V}_c} \sum_{x=1}^N \exp(\mathbf{u}_x^T \mathbf{V}_c) \times \frac{\partial}{\partial \mathbf{V}_c} \sum_{x=1}^N \exp(\mathbf{u}_x^T \mathbf{V}_c)$$

$$\begin{aligned} \therefore \frac{\partial}{\partial v_c} \log p(o|c) &= u_o - \frac{\sum_{x=1}^V \exp(u_x^T v_c) \cdot u_x}{\sum_{w=1}^V \exp(u_w^T v_c)} \\ &= u_o - \underbrace{\frac{\sum_{x=1}^V \frac{\exp(u_x^T v_c)}{\sum_{w=1}^V \exp(u_w^T v_c)} u_x}_{\text{softmax function}} \\ &= u_o - \underbrace{\sum_{x=1}^V p(x|c) u_x}_{\text{expectation } E[x] = p(x) \cdot x} \\ &= \text{observed} - \text{expected} \quad (\text{e.g., } y - y_{\text{pred}}) \end{aligned}$$