

자연어: 모든 시.화.예시 자연이 만들어낸 언어

NLP: 자연어를 컴퓨터가 이해하고 그 의미를 이해하는 기술

구체. 기본 접근

if (나쁜 패턴) ->

규칙적으로 검색/분석

특수 목적 접근

TF-IDF 키워드 추출

인공 지능을 위한 선택 방법

자연어 처리

- 개행문자 제거
- 특수문자 제거
- 공백 제거
- 중복 표현 제거 (ㅋㅋㅋㅋ)
- 이메일 링크 제거
- 제목 제거
- 불용어 제거
- 조사 제거
- 띄어쓰기, 문장 분리 보정
- 사전 구축

Tokenizing

- 글자 tokenizing
- 형태소 tokenizing
- n-gram tokenizing
- word piece tokenizing

Lexical analysis

- 어휘 분석
- 품사 태소 분석
- 개체명 인식
- 상호 참조

Syntactic analysis

- 구문 분석

Semantic analysis

- 의미 분석

one hot encoding

단어 벡터가 sparse 해서 단어를 가지는 의미를 벡터공간에 표현하기

word 2 vec

- 단어 벡터 3차원 공간을 띄우면, 그 단어의 의미를 파악
- 자연어 의미를 벡터 공간에 임베딩
- OOV (out of vocabulary)에 적용 불가능

Sparse representation

- one hot encoding
- 차원의 저주

Dense representation

- word embedding
- 한정된 차원으로 표현 가능
- 의미 관계 파악 가능



단어의 의미가 벡터를 표현함으로써 비슷한 연산 가능

word Embedding 기초

word similarity (사람이 하기) ^{비교} (→) 코딩 가능

analogy test 벡터 연산을 통해 가능

Gensim word2vec library

FastText

- 단어를 n-gram으로 나누어 학습
- 단어를 n-gram으로 분리한 후, 모든 n-gram vector의 평균을 통해 단어 벡터를 찾음

orange \rightarrow <orange> \rightarrow

2-gram	3-gram	4-gram	5-gram	Origin
<O, or, ra, an, ng, we, e>	<or, oru, ran, ang, nge, we>	<ora, aran, rang, ange, nge>	<oraa, orang, range, ange>	orange

단어 벡터와 단어의
유사도를 구함

total n-gram vocabsize 22

그래서 oranges, orange는 orange와 유사한 벡터가 생성

오답률, OOV, 등가 항목이 적은 단어에 적용

word2vec < fast text

Word Embedding의 한계

단어의 학습이 어려움