

정규 교육 세미나

ToBig's 10기 정윤희

Nature Language Process Basic

Contents

Unit 01 | NLP Overview

Unit 02 | Process

Unit 03 | Tokenizing

Unit 04 | Embedding

Unit 05 | Similarity

Unit 06 | Assignment

Unit 01 | NLP Overview



Unit 01 | NLP Overview



NLP 기반 기술 이용 !



Unit 01 | NLP Overview



NLP가 대체 뭘까?



Unit 01 | NLP Overview

Natural Language Processing ?

자연어 처리 또는 자연 언어 처리는
인간의 언어 현상을 컴퓨터와 같은 기계를 이용해서 모사 할수 있도록 연구하고
이를 구현하는 인공지능의 주요 분야 중 하나 !

Unit 01 | NLP Overview

Natural Language Processing ?

NLP가 핫한 이유는?

자연어 처리 또는 자연 언어 처리
인간의 언어 현상을 컴퓨터와 같은 기계를 이용해서 모사 할수 있도록 연구하고
이를 구현하는 인공지능의 주요 분야 중 하나 !

Unit 01 | NLP Overview



- 자연어 이해 및 자연어 처리는 인공지능 분야에 있어서 필수적

Natural Language Processing

지식 기반

통계 기반

딥러닝 기반

사례 기반

머신러닝 기반

Unit 01 | NLP Overview



Contents

Unit 01 | NLP Overview

Unit 02 | Process

Unit 03 | Tokenizing

Unit 04 | Embedding

Unit 05 | Similarity

Unit 06 | Assignment

Unit 02 | Process

NLP는 대체적으로 어떠한 과정을 거칠까요?

Data Collection

Tokenizing

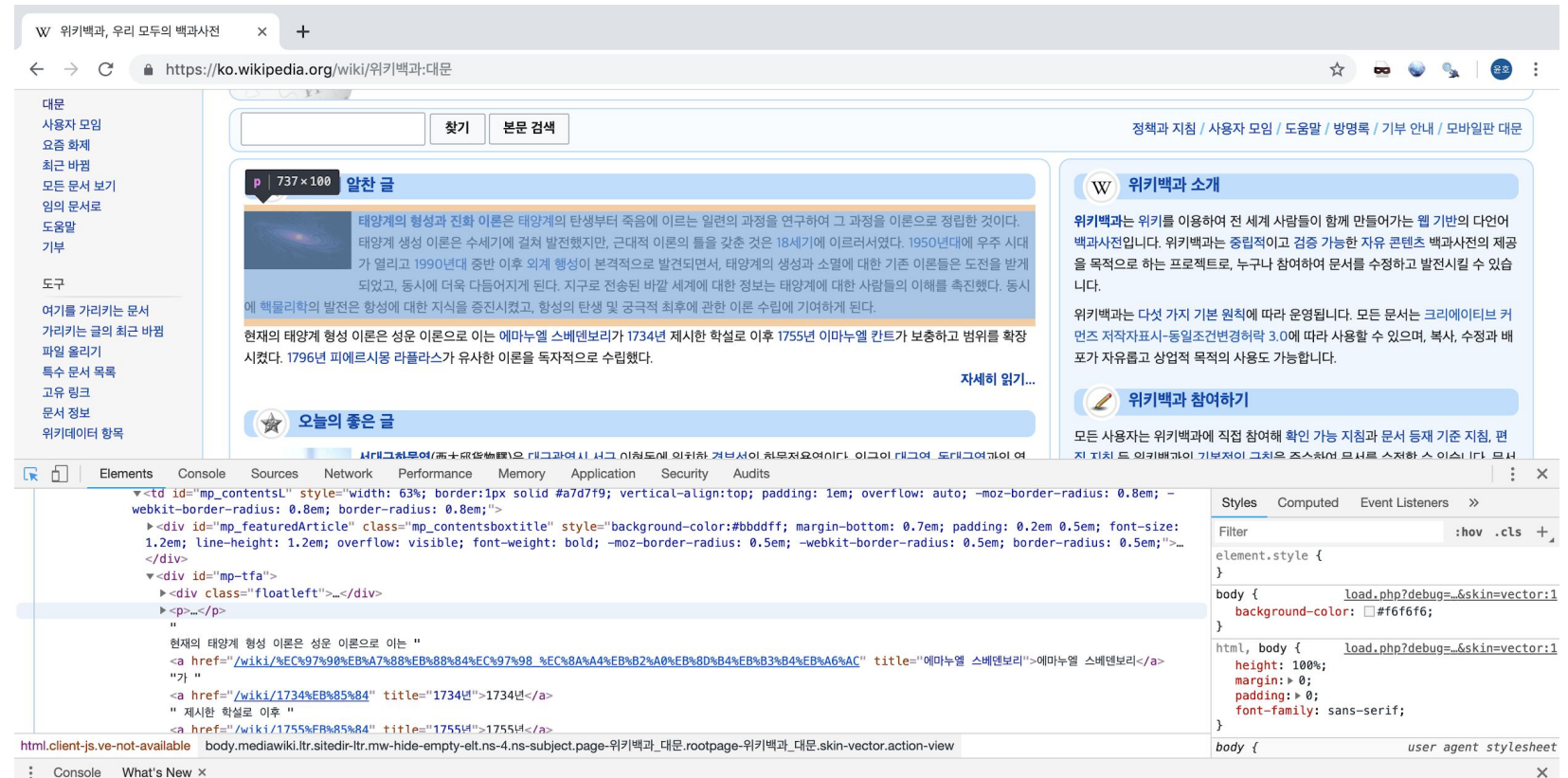
Embedding

Similarity

Network

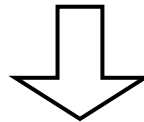
Unit 02 | Process

Data Collection



Tokenizing

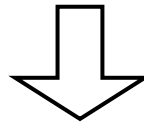
투빅스는 과제지옥입니다.



‘투빅스’, ‘는’, ‘과제’, ‘지옥’, ‘입니’, ‘다.’

Embedding

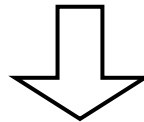
투빅스는 과제지옥입니다.



‘투빅스’ : [0.1234, 0.1234] ‘는’ : [0.5678, 0.1234] ‘과제’ : [0.9012, 0.4321]
‘지옥’ : [0.3456, 0.1764] ‘입니’ : [0.7890, 0.1567] ‘다.’ : [0.1234, 0.3999]

Similarity

‘투빅스’ : [0.1234, 0.1234] ‘는’ : [0.5678, 0.1234] ‘과제’ : [0.9012, 0.4321]
‘지옥’ : [0.3456, 0.1764] ‘입니’ : [0.7890, 0.1567] ‘다.’ : [0.1234, 0.3999]

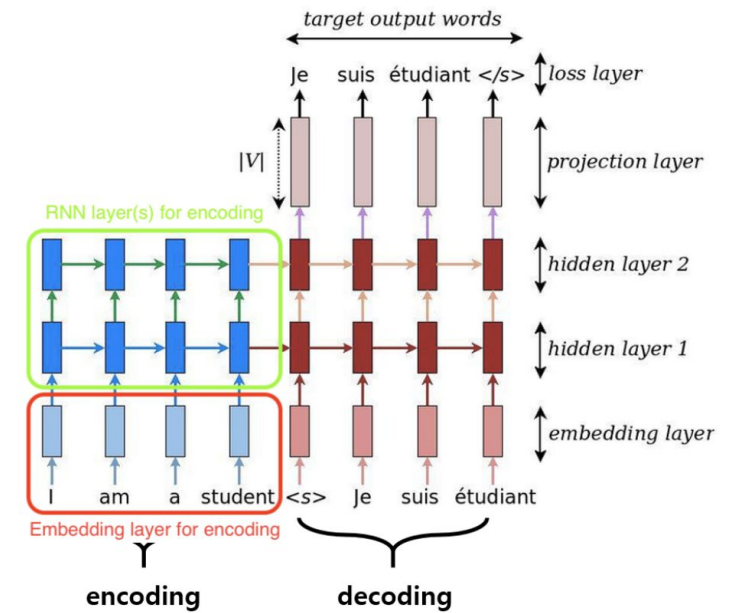
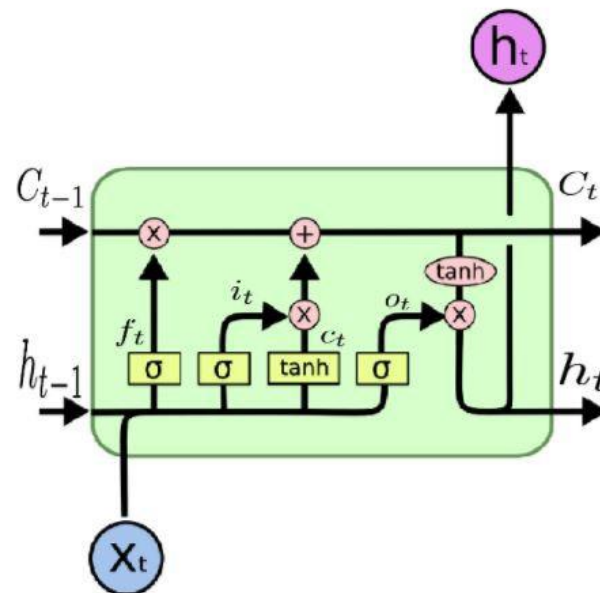
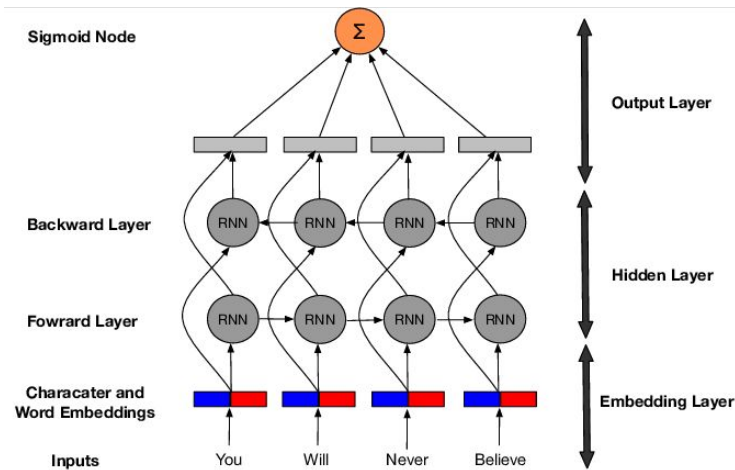


‘과제’ : [0.9012, 0.4321], ‘지옥’ : [0.3456, 0.1764]

코사인 유사도에 따르면, 이 두 단어는 유사하다고 판단할 수 있다.

Unit 02 | Process

Network



Contents

Unit 01 | NLP Overview

Unit 02 | Process

Unit 03 | Tokenizing

Unit 04 | Embedding

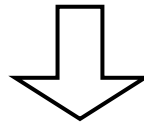
Unit 05 | Similarity

Unit 06 | Assignment

Unit 03 | Tokenizing

Tokenizing

투빅스는 과제지옥입니다.



‘투빅스’, ‘는’, ‘과제’, ‘지옥’, ‘입니’, ‘다.’

Unit 03 | Tokenizing

Tokenizing

투빅스는 과제지옥입니다

특정 기준에 의해서
Text → Token

‘투빅스’, ‘는’, ‘과제’, ‘지옥’, ‘입니’, ‘다.’

Unit 03 | Tokenizing

Tokenizing

투빅스는 과제지옥입니다.

자소/음소, 형태소, 단어, 문장, 문서... etc



‘투빅스’, ‘는’, ‘과제’, ‘지옥’, ‘입니’, ‘다.’

Unit 03 | Tokenizing

English

NLTK

Korean

KONLPY

Unit 03 | Tokenizing

Kkma

Twitter(Okt)

Komoran

Hannanum

Mecab

Morphs

Nouns

Pos Tagging

Unit 03 | Tokenizing

아버지가방에들어가신다.

Hannanum	Kkma	Komoran	Mecab	Twitter
아버지가방에 들어가 / N	아버지 / NNG	아버지가방에 들어가신다 / NNP	아버지 / NNG	아버지 / Noun
이 / J	가방 / NNG		가 / JKS	가방 / Noun
시ㄴ다 / E	에 / JKM		방 / NNG	에 / Josa
	들어가 / VV		에 / JKB	들어가신 / Verb
	시 / EPH		들어가 / VV	다 / Eomi
	ㄴ다 / EFN		신다 / EP+EC	

Unit 03 | Tokenizing

아버지가방에들어가신다.

Hannanum	Kkma	Komoran	Mecab	Twitter
아버지가방에 들어가 / N	아버지 / NNG	아버지가방에 들어가신다 / EP	아버지 / NNG	아버지 / Noun
이 / J	가방 / NNG		가 / JKS	가방 / Noun
시ㄴ다 / E	에 / JKM		방 / NNG	에 / Josa
	들어가 / VV		에 / JKB	들어가신 / Verb
	시 / EPH		들어가 / VV	다 / Eomi
	ㄴ다 / EFN		신다 / EP+EC	

실습 코드 !!

Unit 03 | Tokenizing

딥러닝 기반 형태소 분석기

kakao Tech 블로그 오픈소스 오픈API 기술행사 Search... 🔍

< Back to Posts 📱 🌐 🐦 🍏

kakao의 오픈소스 Ep9 - Khaiii : 카카오의 딥러닝
기반 형태소 분석기

opensource , khaiii , deep-learning , cnn , ai

“카카오의 오픈소스를 소개합니다” 아홉 번째는 jamie.lim과 자연어 처리 파트 동료들이 함께 개발한 khaiii(Kakao Hangul Analyzer III)입니다.

khaiii는 세종 코퍼스를 이용하여 CNN(Convolutional Neural Network, 합성곱 신경망) 기술을 적용해 학습한 형태소 분석기입니다. 디코더를 C++로 구현하여 GPU 없이도 비교적 빠르게 동작하며, Python 바인딩을 제공하고 있어서 편리하게 이용하실 수 있습니다.

Fork me on GitHub

<https://github.com/kakao/khaiii>

Contents

Unit 01 | NLP Overview

Unit 02 | Process

Unit 03 | Tokenizing

Unit 04 | Embedding

Unit 05 | Similarity

Unit 06 | Assignment

Unit 04 | Embedding

Tokenizing 왜 했지?

Unit 04 | Embedding

자연어 처리를 위한 의미단위를 만들기 위해 !

Unit 04 | Embedding

그런데 컴퓨터가 인간의 언어를 어떻게 이해할 수 있을까?

Unit 04 | Embedding

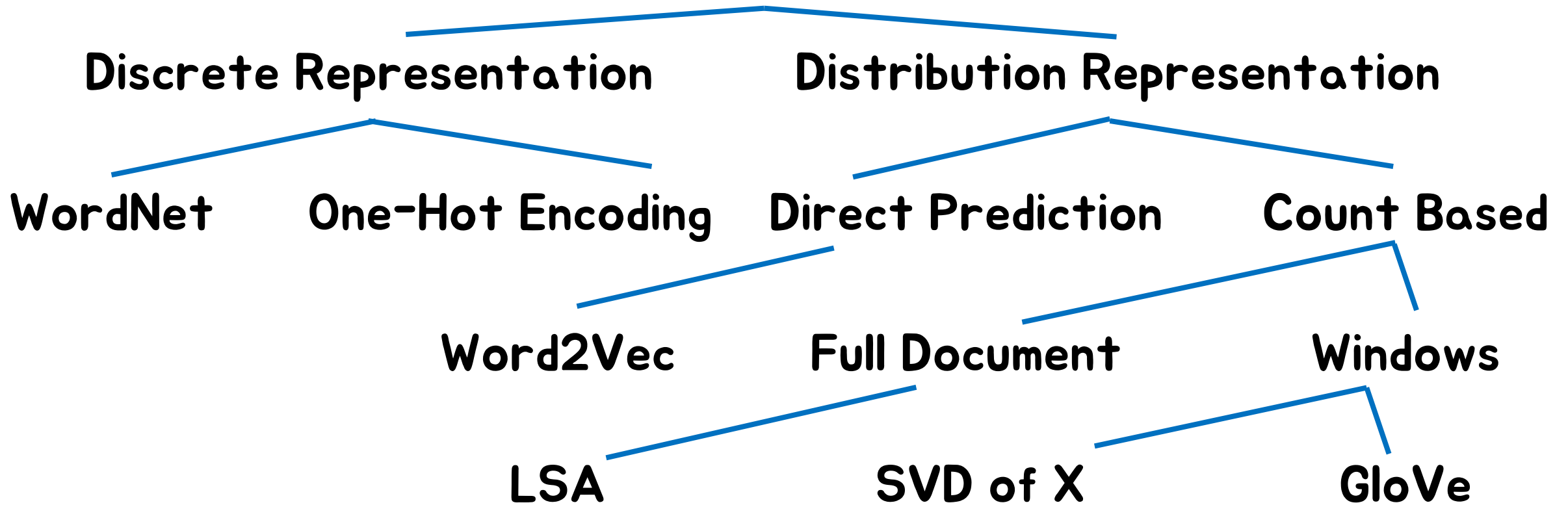
컴퓨터가 처리할 수 있는 것은 수치분

Unit 04 | Embedding

컴퓨터가 **언어의 특성**을 이해할 수 있도록
각 Token마다 **수치**를 부여

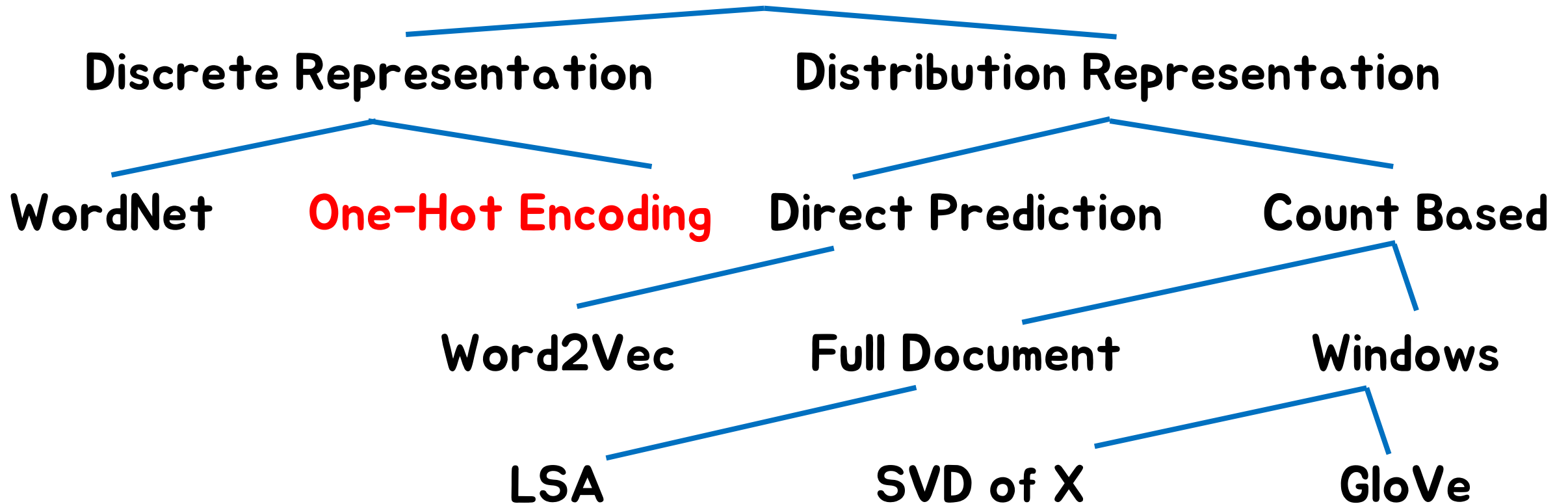
Unit 04 | Embedding

Word Representation



Unit 04 | Embedding

Word Representation



One-Hot Encoding의 문제점

1. n 개 token \rightarrow n 개 feature. 불필요한 계산이 많다
2. 유사도 측정이 어려워 유의어, 반의어 등의 언어적 특성을 고려하기 힘들다.

Unit 04 | Embedding

해질 무렵 바람도 몹시 불던 날 집에 돌아오는 길 버스 창가에 앉아

Unit 04 | Embedding

['해', '질', '무렵', '바람', '도', '몹시', '불', '던', '날', '집', '에', '돌아오는', '길', '버스', '창가', '에', '앉아']

Unit 04 | Embedding

One-Hot Encoding

'해' : [1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0]
 '질' : [0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0]
 '무렵' : [0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0]
 '바람' : [0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0]
 '도' : [0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0]
 '몹시' : [0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0]
 '불' : [0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0]
 '던' : [0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0]
 '날' : [0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0]
 '집' : [0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0]
 '에' : [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0]
 '돌아오는' : [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0]
 '길' : [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0]
 '버스' : [0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0]
 '창가' : [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0]
 '에' : [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0]
 '앞아' : [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1]

한 문장을 표현하기 위해
17개의 feature 필요

해질 무렵 바람도 몹시 불던 날 집에 돌아오는 길 버스 창가에 앉아 불어오는
바람 어찌지도 못한 채 난 그저 멍할 뿐이었지 난 왜 이리 바보인지 어리석은
지 모진 세상이란 걸 아직 모르는 지 터지는 울음 입술 물어 삼키며 내려야지
하고 일어설 때 저 멀리 가까워 오는 정류장 앞에 희미하게 일렁이는
언제부터 기다렸는지 알 수도 없는 발만 동동 구르고 있는 그댈 봤을 때 나는
아무 말도 못하고 그댈 안고서 그냥 눈물만 흘러 자주 눈물이 흘러 이대로
영원히 있을 수만 있다면 오 그대여 그대여서 고마워요 낙엽이 뒹굴고 있는
정류장 앞에 희미하게 일렁이는 까치발 들고 내 얼굴 찾아 헤매는 내가
사준 옷을 또 입고 온 그댈 봤을 때 나는 아무 말도 못하고 그댈 안고서 그냥
눈물만 흘러 자주 눈물이 흘러 이대로 영원히 있을 수만 있다면 오 그대여
그대여서 고마워요

Unit 04 | Embedding

One-Hot Encoding

'해' : [1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0]
 '질' : [0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0]
 '무렵' : [0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0]
 '바람' : [0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0]
 '도' : [0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0]
 '몹시' : [0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0]
 '불' : [0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0]
 '던' : [0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0]
 '날' : [0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0]
 '집' : [0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0]
 '에' : [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0]
 '돌아오는' : [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0]
 '길' : [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0]
 '버스' : [0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0]
 '창가' : [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0]
 '에' : [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0]
 '앞아' : [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1]

Discrete Representation
Sparse Vector
Inner Product = 0
(Independent)

Unit 04 | Embedding

One-Hot Encoding

'해' : [1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0]
 '질' : [0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0]
 '무렵' : [0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0]
 '바람' : [0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0]
 '도' : [0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0]
 '몹시' : [0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0]
 '불' : [0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0]
 '던' : [0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0]
 '날' : [0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0]
 '집' : [0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0]
 '에' : [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0]
 '돌아오는' : [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0]
 '길' : [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0]
 '버스' : [0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0]
 '창가' : [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0]
 '에' : [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0]
 '앞아' : [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1]

- 차원이 너무 커지고 불필요한 계산이 많아짐
- 유사도 측정이 어려워 유의어, 반의어 등의 언어적 특성을 고려하지 못함

Unit 04 | Embedding

효과적인 방법이 없을까?

Unit 04 | Embedding

Word Embedding이란 ?

단어를 좀 더 **조밀한 차원**에 **벡터**로 표현하는 것

Unit 04 | Embedding

Word Representation

Discrete Representation

WordNet

One-Hot Encoding

Distribution Representation

Direct Prediction

Count Based

Word2Vec

Full Document

Windows

LSA

SVD of X

GloVe

Unit 04 | Embedding

GloVe

BERT

Word2Vec

FastText

Unit 04 | Embedding

GloVe

BERT

Word2Vec

FastText

Unit 04 | Embedding

Word2Vec : 말 그대로 Word to Vector

Unit 04 | Embedding

How ?

Word2Vec : 말 그대로 Word to Vector

Unit 04 | Embedding

1. Word & Neighbor (CBOW, Skip Gram)

2. Neural Network

Unit 04 | Embedding

CBOW

Skip Gram

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

EX) Skip Gram

Unit 04 | Embedding

Window Size

ex) window size = 2

원수는 외나무다리에서 만난다.

Unit 04 | Embedding

Window Size

‘원수’, ‘는’, ‘외나무다리’
‘에서’, ‘만난다’

Center Word	Neighbor Words
‘원수’	‘는’, ‘외나무다리’
‘는’	‘원수’, ‘는’, ‘외나무다리’
‘외나무다리’	‘원수’, ‘는’, ‘에서’, ‘만난다’
‘에서’	‘는’, ‘외나무다리’, ‘만난다’
‘만난다’	‘외나무다리’, ‘에서’

Unit 04 | Embedding

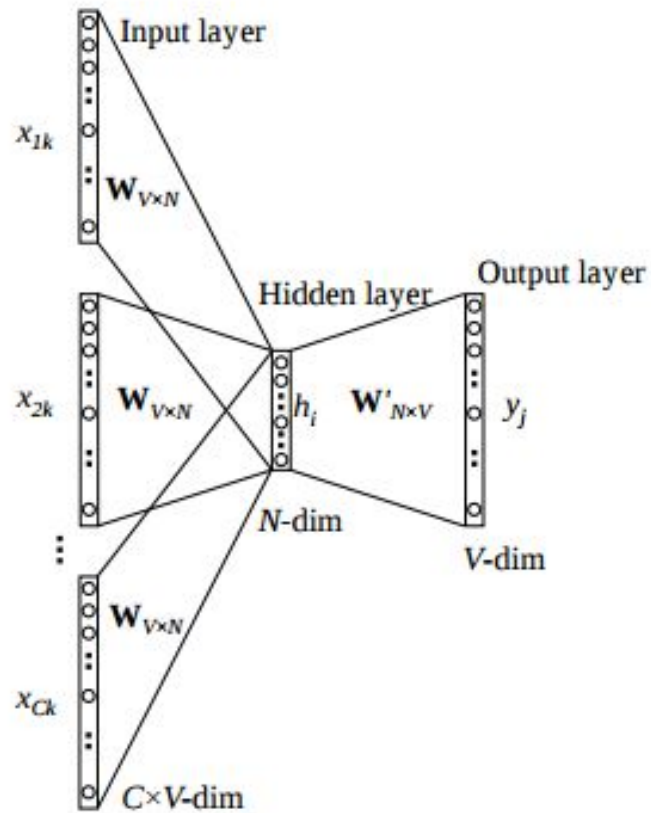
CBOW

원수는 _____에서 만난다.

Unit 04 | Embedding

CBOW

Input : Neighbor Words



Target : Center Word

Unit 04 | Embedding

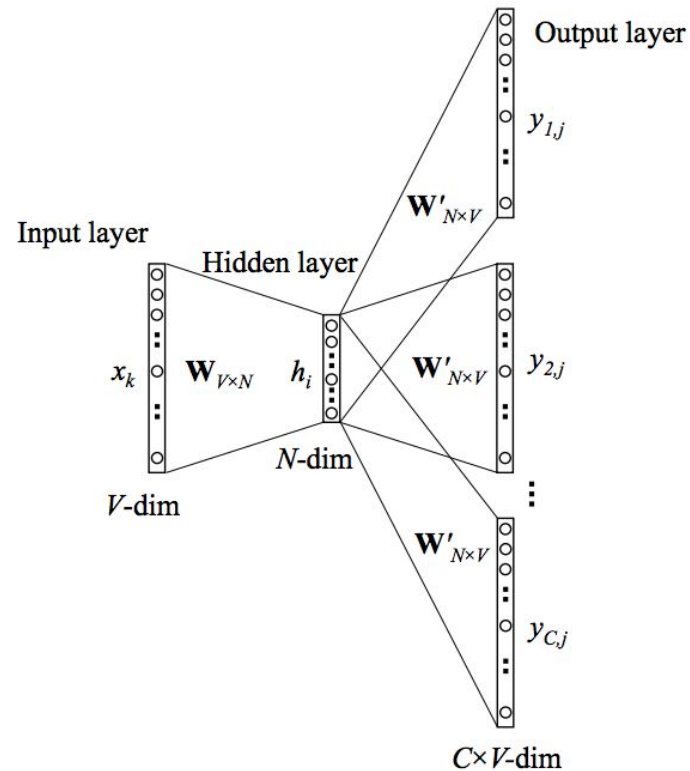
Skip Gram

___외나무다리___.

Unit 04 | Embedding

Skip Gram

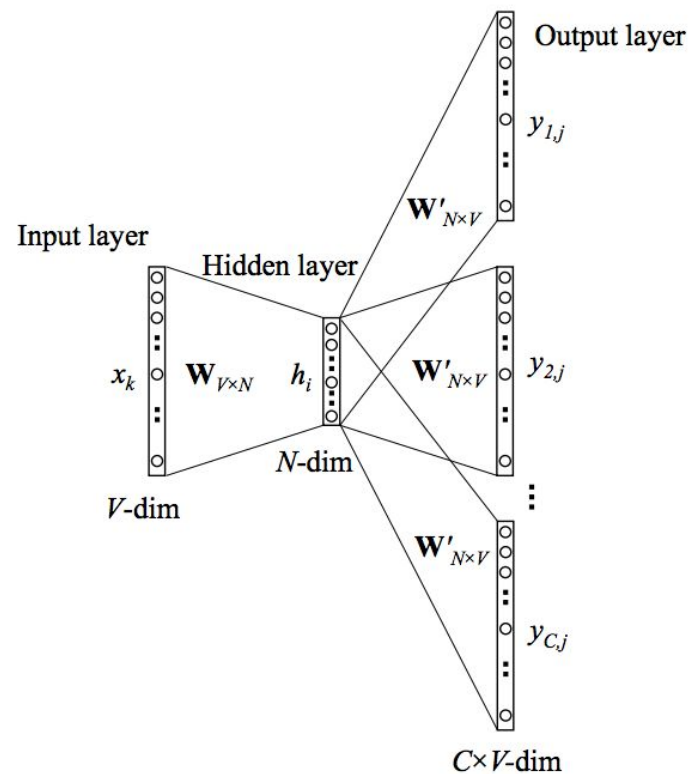
Input : Center Word



Target : Neighbor Words

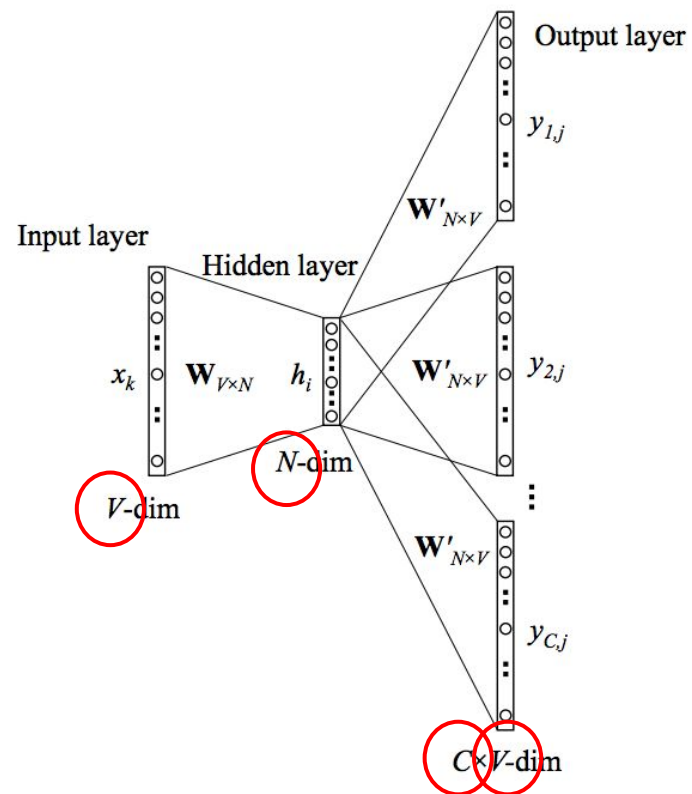
Unit 04 | Embedding

Neural Network



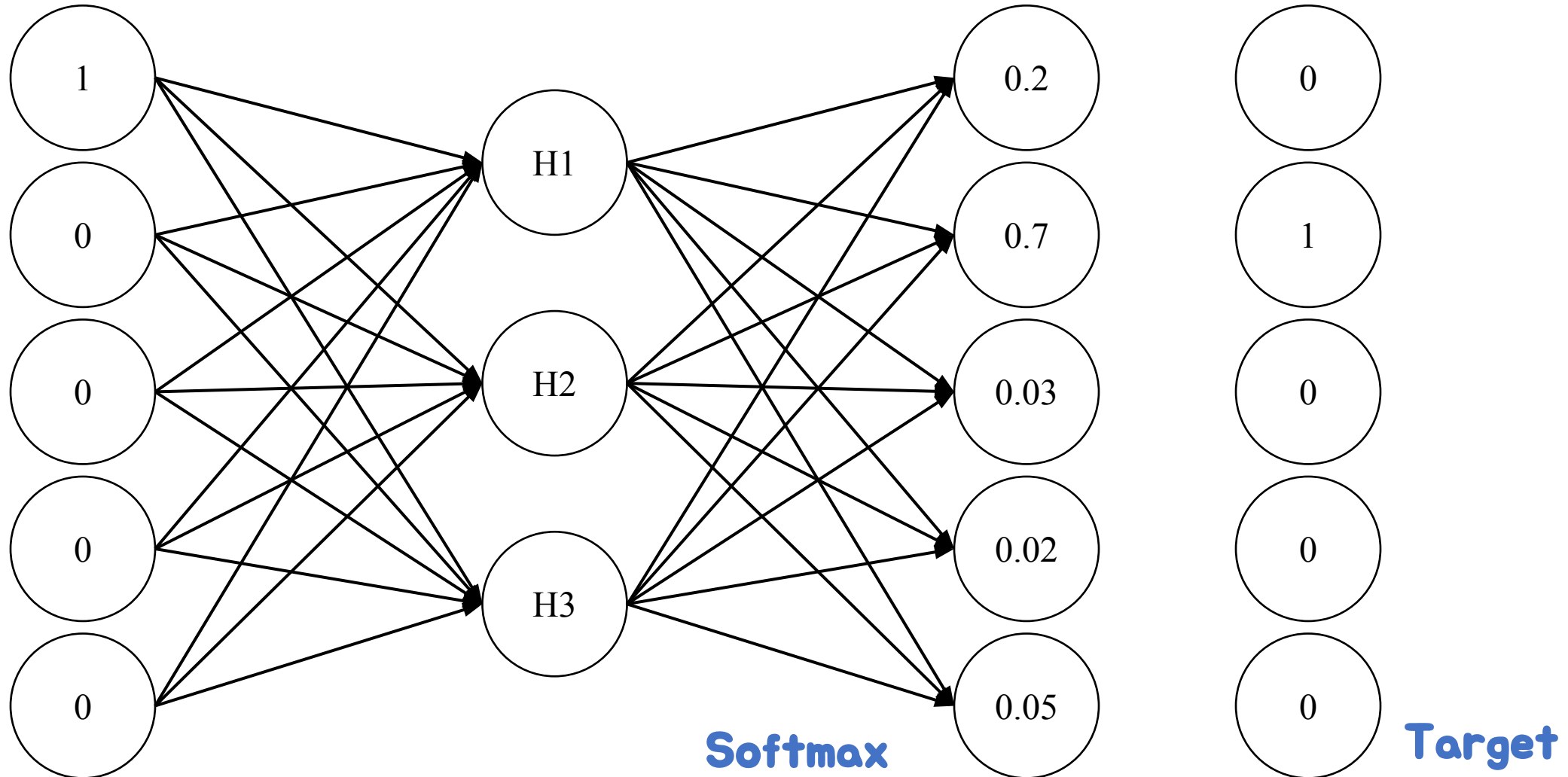
Unit 04 | Embedding

Neural Network



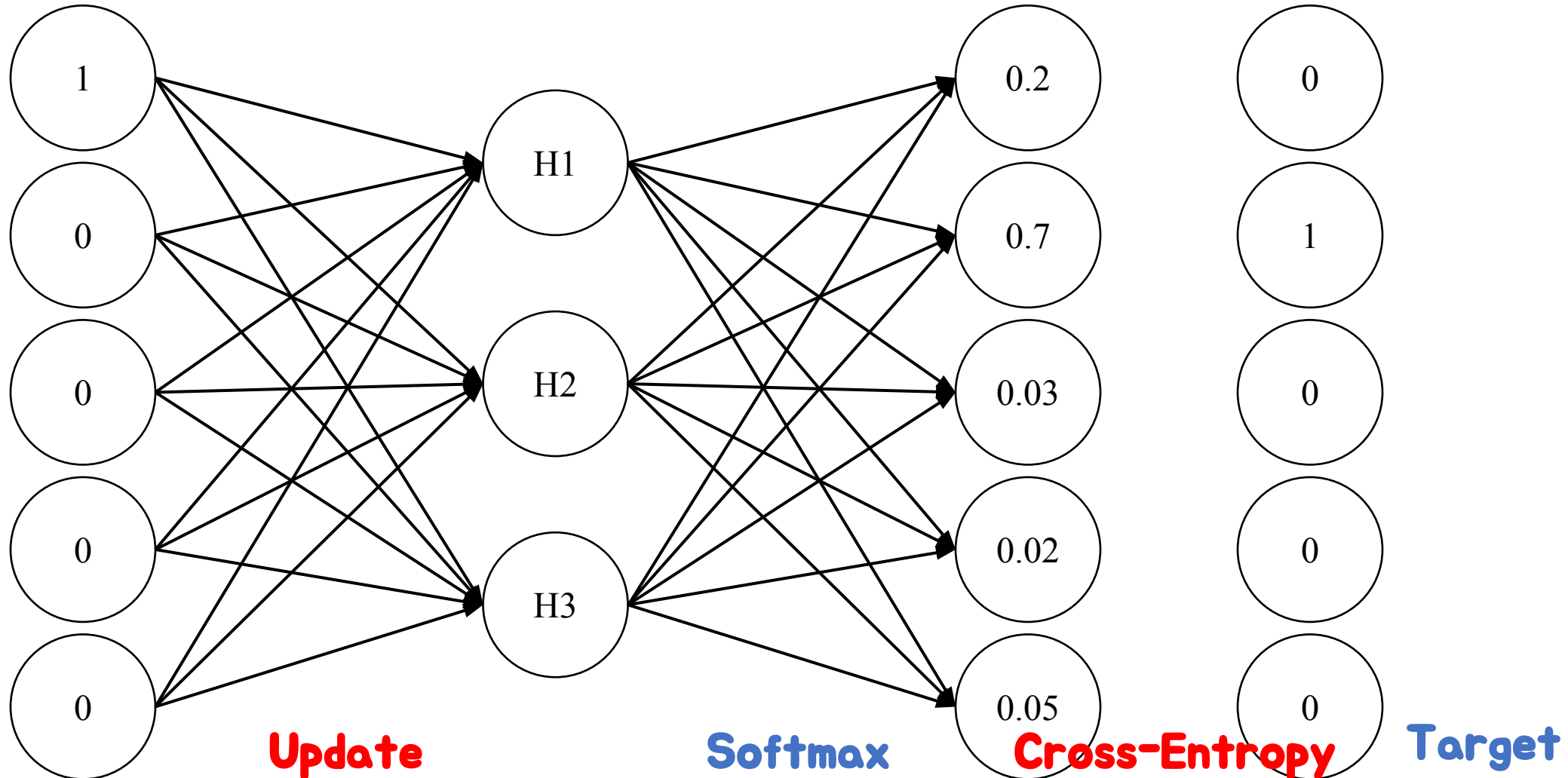
Unit 04 | Embedding

**One-Hot
Encoding**



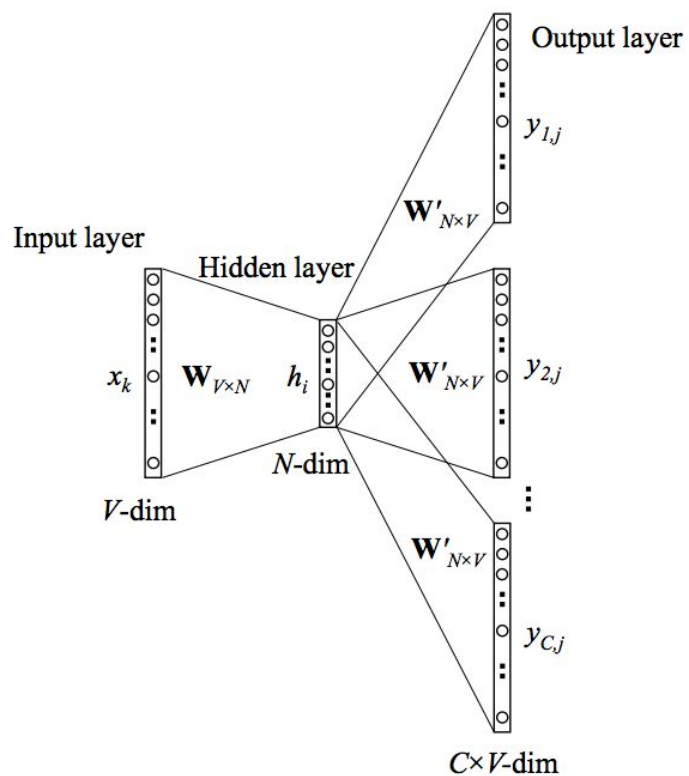
Unit 04 | Embedding

One-Hot
Encoding



Unit 04 | Embedding

Word Vector

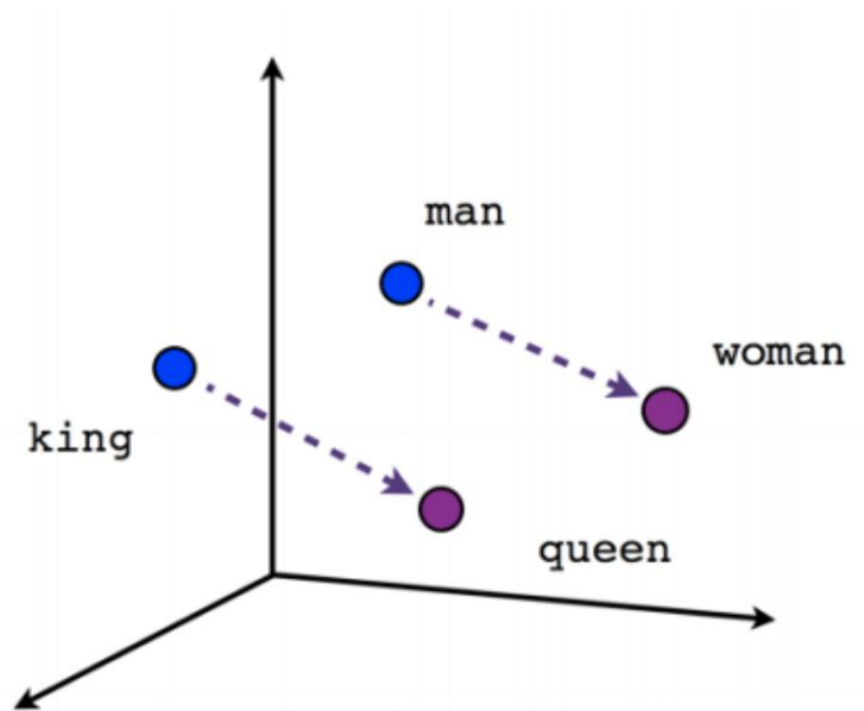


H1

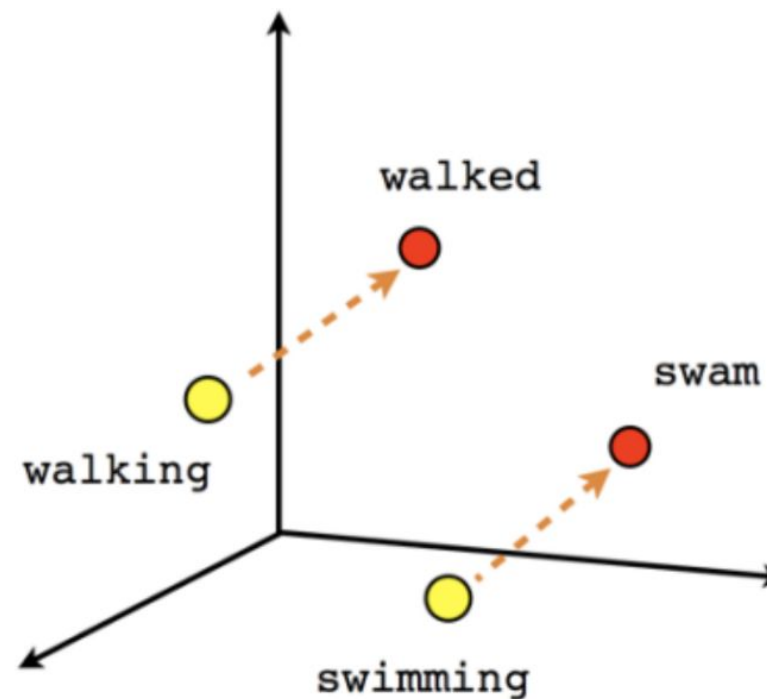
H2

H3

Unit 04 | Embedding



Male-Female



Verb tense

Unit 04 | Embedding

Word2Vec의 문제점 ?

1. 한번에 하나의 동시 출현만 계산해서 업데이트 -> 전체적인 통계(빈도 수 등) 정보 이용x -> 비효율적인 면이 있지 않은가?
2. train corpus에 존재하지 않았던 단어의 벡터를 만들어낼 수 없다.

Unit 04 | Embedding

GloVe

전체적인 **통계 정보**를 이용해보자 !

학습 말뭉치에서 동시에 **같이 등장한 단어의 빈도**를 각각 세어서
전체 말뭉치의 단어 개수로 나눠준 **동시등장확률**을 고려하자 !

Unit 04 | Embedding

GloVe

- Example corpus:
 - I like deep learning.
 - I like NLP.
 - I enjoy flying.

counts	I	like	enjoy	deep	learning	NLP	flying	.
I	0	2	1	0	0	0	0	0
like	2	0	0	1	0	1	0	0
enjoy	1	0	0	0	0	0	1	0
deep	0	1	0	0	1	0	0	0
learning	0	0	0	1	0	0	0	1
NLP	0	1	0	0	0	0	0	1
flying	0	0	1	0	0	0	0	1
.	0	0	0	0	1	1	1	0

Probability and Ratio	$k = solid$	$k = gas$	$k = water$	$k = fashion$
$P(k ice)$	1.9×10^{-4}	6.6×10^{-5}	3.0×10^{-3}	1.7×10^{-5}
$P(k steam)$	2.2×10^{-5}	7.8×10^{-4}	2.2×10^{-3}	1.8×10^{-5}
$P(k ice)/P(k steam)$	8.9	8.5×10^{-2}	1.36	0.96

FastText

단어가 아닌 **단어내부의 n-gram**을 최소단위로 !

Unit 04 | Embedding

FastText

'apple'

min = 3, max = 6

'<ap', 'app', 'appl', 'apple', 'apple>', 'ppl',
'pple', 'pple>', 'ple', 'ple>', 'le>'

Unit 04 | Embedding

FastText

1. train corpus에 존재하지 않았던 단어의 벡터를 만들어낼 수 있다.(word2vec,glove에서는 불가능)
2. 희소한 단어에 대해 더 좋은 word embedding이 가능하다.

Unit 04 | Embedding

BERT

google-research / bert

Watch 726 Star 12,741 Fork 2,727

Code Issues 191 Pull requests 21 Projects 0 Wiki Insights

TensorFlow code and pre-trained models for BERT <https://arxiv.org/abs/1810.04805>

nlp google natural-language-processing natural-language-understanding tensorflow

103 commits 1 branch 0 releases 26 contributors Apache-2.0

Branch: master New pull request Create new file Upload files Find file Clone or download

jacobdevlin-google Merge pull request #431 from dalequark/colab-tfhub Latest commit ffbd2a 21 days ago

.gitignore	Initial BERT release	4 months ago
CONTRIBUTING.md	Initial BERT release	4 months ago
LICENSE	Initial BERT release	4 months ago
README.md	Changed colab link	21 days ago
init.py	Initial BERT release	4 months ago
create_pretraining_data.py	Adding TF Hub support	26 days ago
extract_features.py	Running through pyformat to meet Google code standards	4 months ago
modeling.py	Adding TF Hub support	26 days ago
modeling_test.py	Adding SQuAD 2.0 support	4 months ago
multilingual.md	Padding examples for TPU eval/predictions and checking case match	3 months ago
optimization.py	Padding examples for TPU eval/predictions and checking case match	3 months ago
optimization_test.py	Initial BERT release	4 months ago

<https://github.com/google-research/bert>

Contents

Unit 01 | NLP Overview

Unit 02 | Process

Unit 03 | Tokenizing

Unit 04 | Embedding

Unit 05 | Similarity

Unit 06 | Assignment

Similarity Analysis

Distance

**Cosine
Similarity**

Unit 05 | Similarity

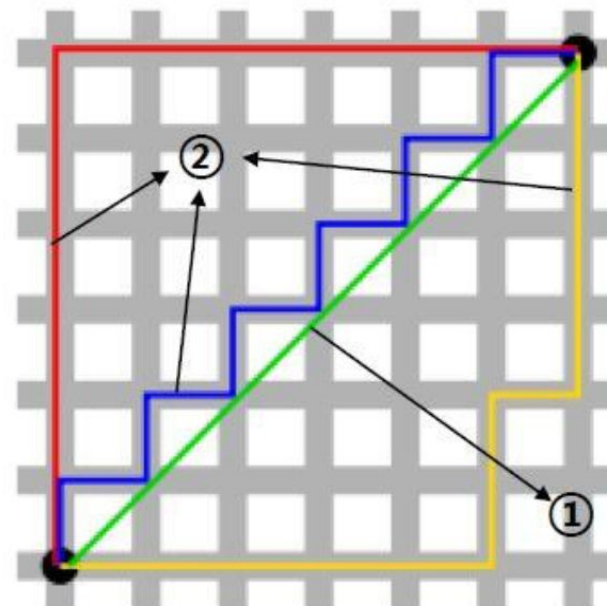
Distance

$$d(x, y) = (\sum_{i=1}^p (x_i - y_i)^2)^{1/2}$$

$$d(x, y) = \sum_{i=1}^p |x_i - y_i|$$

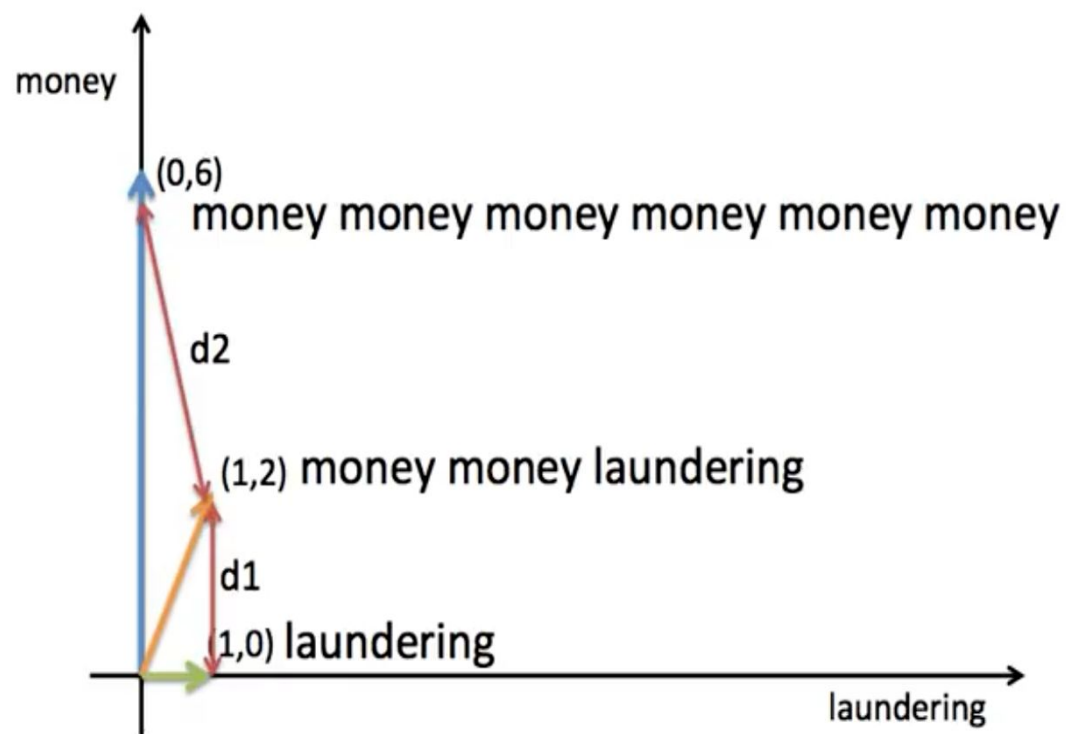
$$d(x, y) = (\sum_{i=1}^p (x_i - y_i)^2 / s_i^2)^{1/2}$$

$$d(x, y) = (\sum_{i=1}^p (x_i - y_i)^m)^{1/m}$$



Unit 05 | Similarity

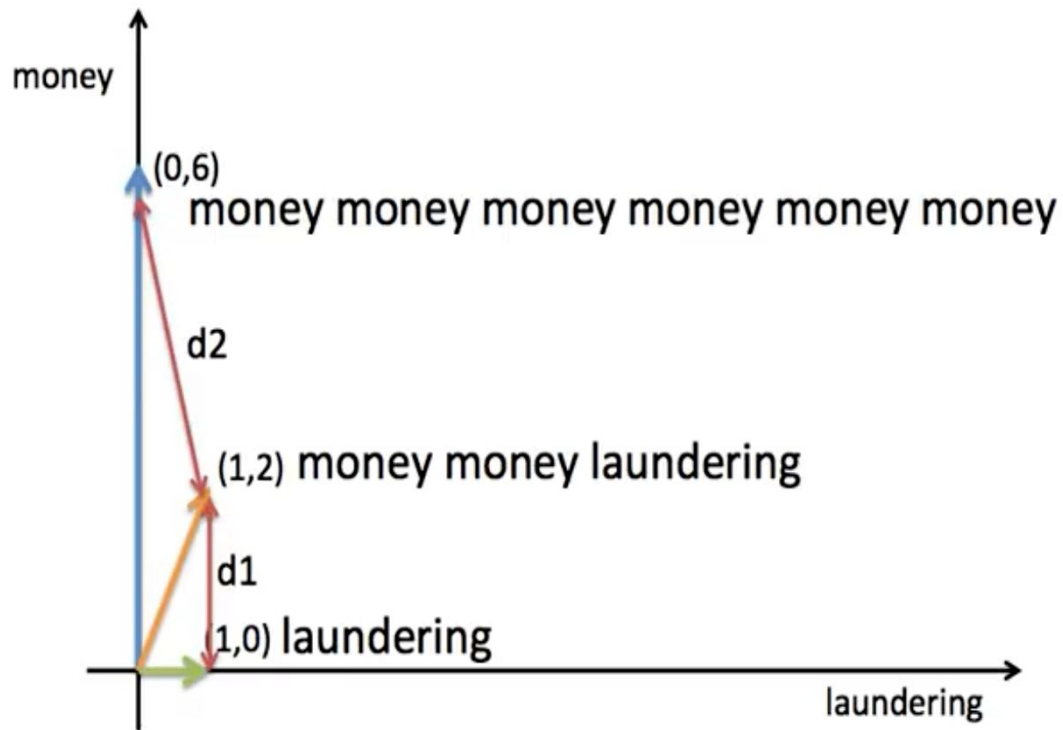
Cosine Similarity



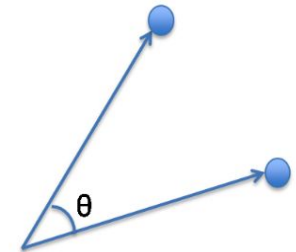
$d2 > d1$
하지만 money money laundrying은
money money money money money와
더 유사하지 않을까?

Unit 05 | Similarity

Cosine Similarity



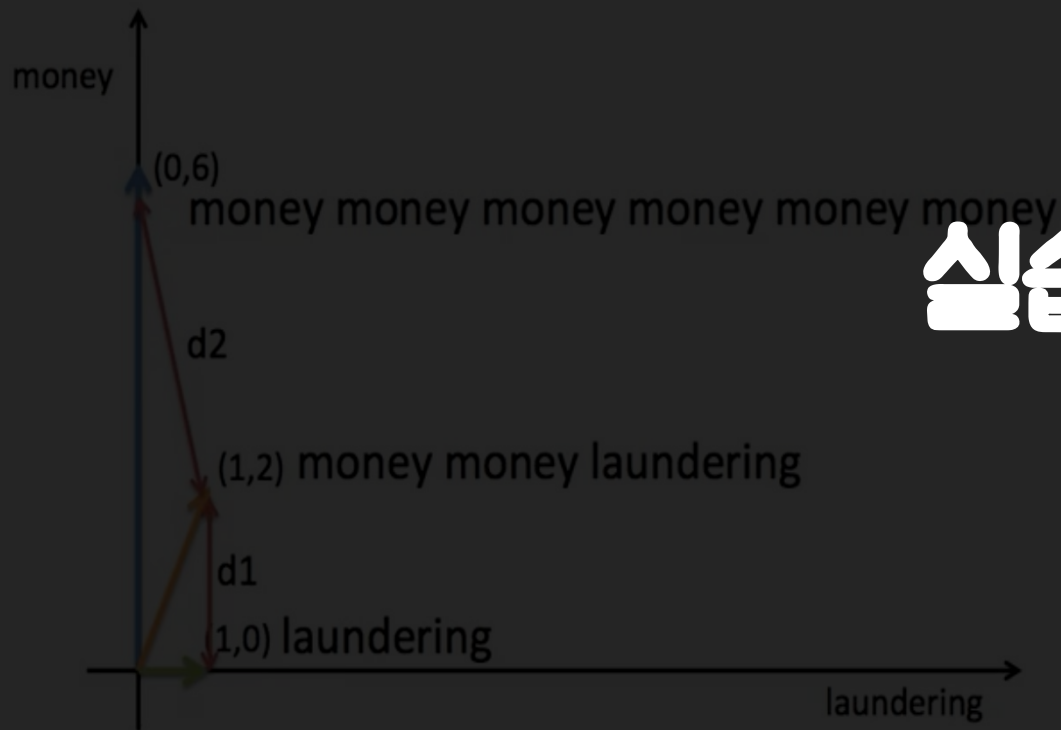
$$\text{sim}(A, B) = \cos(\theta) = \frac{A \cdot B}{\|A\| \|B\|}$$



**벡터의 크기를 무시하고 방향성만 고려하여 계산
두 벡터가 이루고 있는 각도를 활용하여
작을수록 유사하다고 판단**

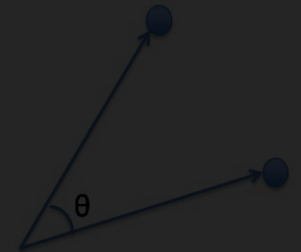
Unit 05 | Similarity

Cosine Similarity



실습 코드 !!

$$\cos(A, B) = \cos(\theta) = \frac{A \cdot B}{\|A\| \|B\|}$$



벡터의 크기를 무시하고 방향성만 고려하여 계산
두 벡터가 이루고 있는 각도를 활용하여,
작을수록 유사하다고 판단

Contents

Unit 01 | NLP Overview

Unit 02 | Process

Unit 03 | Tokenizing

Unit 04 | Embedding

Unit 05 | Similarity

Unit 06 | Assignment

Unit 06 | Assignment

〈 과제1 〉

Word2Vec을 구현한 TensorFlow 코드에서, # 부분에 있는 코드의 의미에 대해 주석달아보기 ! 그리고 하이퍼 파라미터 조정해서 코드 다시 실행해보기 !

〈 과제 2〉

프로젝트 주제 혹은 관심분야를 정해서 해당 주제에 대한 텍스트를 크롤링하고, 전처리와 임베딩을 하고 난 후 그래프 그리기 !(우수과제는 크롤링이나 전처리 과정에서의 센스, 혹은 임베딩 결과를 어떻게 이용하고 분석하느냐에 달려있음. 이 부분은 선택적)

Unit 06 | Assignment

〈Reference〉

- [http://www.datamarket.kr/xe/\(투빅스 9기 안상준 NLP Basic\)](http://www.datamarket.kr/xe/(투빅스%209기%20안상준%20NLP%20Basic))
- <https://www.youtube.com/watch?v=uZ2GtEe-50E>
- <https://ratsgo.github.io/blog/categories/>
- <http://web.stanford.edu/class/cs224n/>
- <https://www.youtube.com/watch?v=sY4YyacSsLc>
- <https://wikidocs.net/22660>
- <https://datascienceschool.net/view-notebook/6927b0906f884a67b0da9310d3a581ee/>

Unit 06 | Assignment

〈IMG〉

- http://techm.kr/bbs/board.php?bo_table=article&wr_id=4051
- <http://www.ndsl.kr/ndsl/issueNdsl/detail.do?techSq=50>
- <https://towardsdatascience.com/seq2seq-model-in-tensorflow-ec0c557e560f>
- <https://hackernoon.com/understanding-architecture-of-lstm-cell-from-scratch-with-code-8da40f0b71f4>
- https://www.researchgate.net/figure/BiDirectional-RNN-architecture-for-detecting-clickbaits_fig1_311430194
- <https://imgur.com/TupGxMI>
- <https://www.analyticsvidhya.com/blog/2017/06/word-embeddings-count-word2veec/>
- <https://www.tensorflow.org/images/linear-relationships.png>
- <http://web.stanford.edu/class/cs224n/>

Q & A

들어주셔서 감사합니다.