



Text Mining & Word Embedding

22.08.30 / DSL 7기 최명헌

1. Text Mining

- What is Text Mining?
- Crawling
- Cleansing & Tokenization

2. Word Embedding

- word2vec
- CBoW
- Skip-Gram
- Hierarchical Softmax
- Negative Sampling
- Others...

3. Visualization

- PCA
- t-SNE

1. Text Mining



What is Text Mining?

Text + Mining = **글자를 채굴하는 것**
글자를 가져오는 것? = Crawling?

BeautifulSoup

 Selenium

1. Text Mining

What is Text Mining?



Text mining is the process of *deriving high-quality information* from text.

– Wikipedia

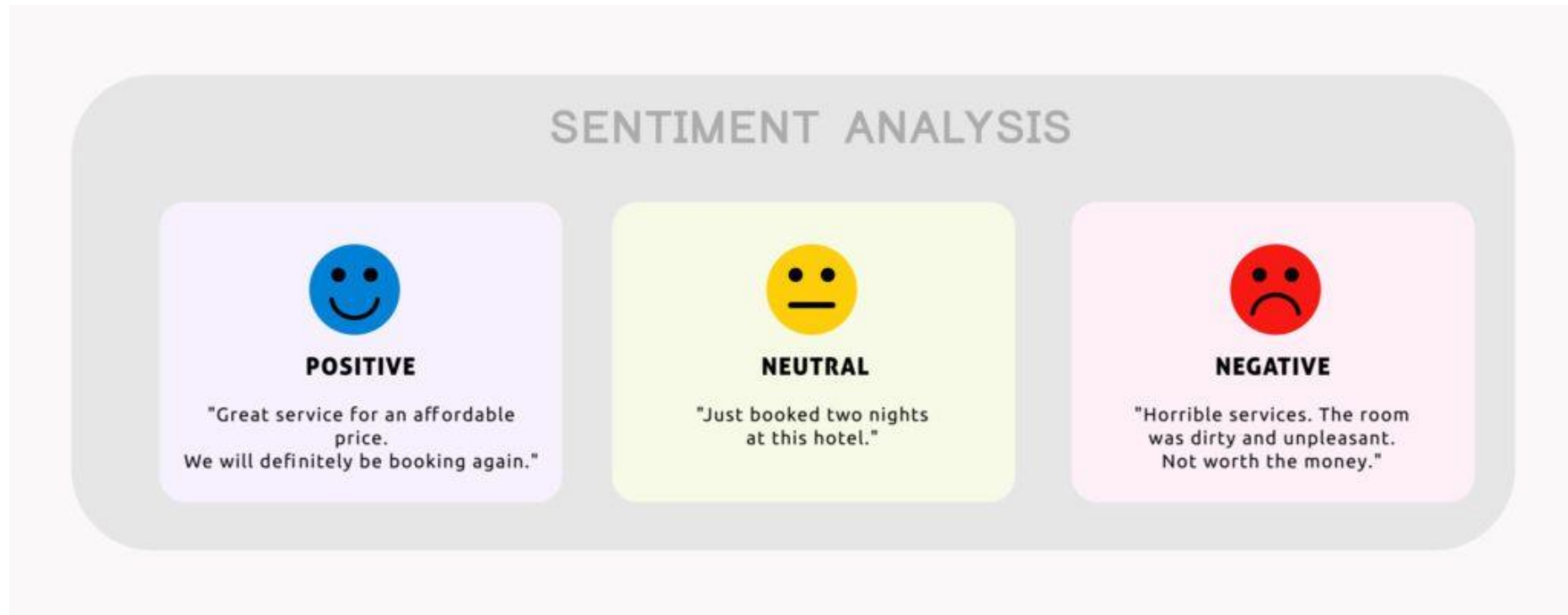


1. Text Mining



What is Text Mining?

High Quality Information 1 : Sentiment Analysis



1. Text Mining



What is Text Mining?

High Quality Information 2 : Topic Analysis

해럴드경제 구독 PICK 해럴드경제 헤드라인

“얼굴사진 어느 게 낫나요?” 갤럭시4 vs 갤럭시22울트라, 비교해보니

입력 2022.08.21. 오후 7:52 기사원문 김현일 기자 >

4 21

본문 요약봇

자동 추출 기술로 요약된 내용입니다. 요약 기술의 특성상 본문의 주요 내용이 제외될 수 있어, 전체 맥락을 이해하기 위해서는 기사 본문 전체보기를 권장합니다.

“얼굴사진 어느 게 낫나요?” 갤럭시4 vs 갤럭시22울트라, 비교해보니

삼성전자의 4세대 폴더블폰이 사전예약이 지난 16일부터 시작된 가운데 카메라 성능에 대한 소비자들의 궁금증을 조금이나마 해소해줄 결과가 나와 눈길을 끈다.

IT팁스터 아이스유니버스는 최근 자신의 SNS에 삼성전자의 '갤럭시 S22 울트라'와 '갤럭시 Z폴드4' 카메라로 동일한 여성의 얼굴을 '인물모드'로 찍은 사진을 올렸다.

삼성전자가 올해 상반기와 하반기에 각각 선보인 최신 '바형 스마트폰'과 200만원대 초고가 '접는 스마트폰' 카메라 간의 대결이어서 이목을 집중시켰다.

“손흥민 믿고 유료화 질렀더니” 육만 엄청 먹다가, 이런일이 8시간전

N 예약

1. Text Mining



What is Text Mining?

High Quality Informations...

Two screenshots of a chat interface with a character named '런디' (Rendi).

Left Screenshot:

- Header: 런어스 챗봇 런디(Learn_D)
- Time: 2:57 am
- User input: 뭐해? (What are you doing?)
- Time: 2:58 am
- User input: 고려대 가고 싶어 (I want to go to Yonsei University)
- Time: 2:57 am
- Character response: 지금은 에너지 충전중이에요! 잘 쉬는 것도 중요하답니다! (I'm currently charging energy! Resting well is also important!)
- Time: 2:57 am

Right Screenshot:

- Header: 런어스 챗봇 런디(Learn_D)
- Time: 2:58 am
- Character response: 런디는 그런 단어 잘 모르겠어요! (Rendi doesn't know that word well!)
- Time: 2:58 am

Two screenshots of a translation interface.

Top Screenshot:

- Language: 한국어 감지 (Korean detected)
- Target language: 영어 (English)
- Text: 데이터 사이언스 랩은 연세대 최고의 학회입니다. (Data Science Lab is Yonsei University's best academic society.)
- Character count: 26 / 5000
- Buttons: 자동완성 (Auto-complete), 번역하기 (Translate)

Bottom Screenshot:

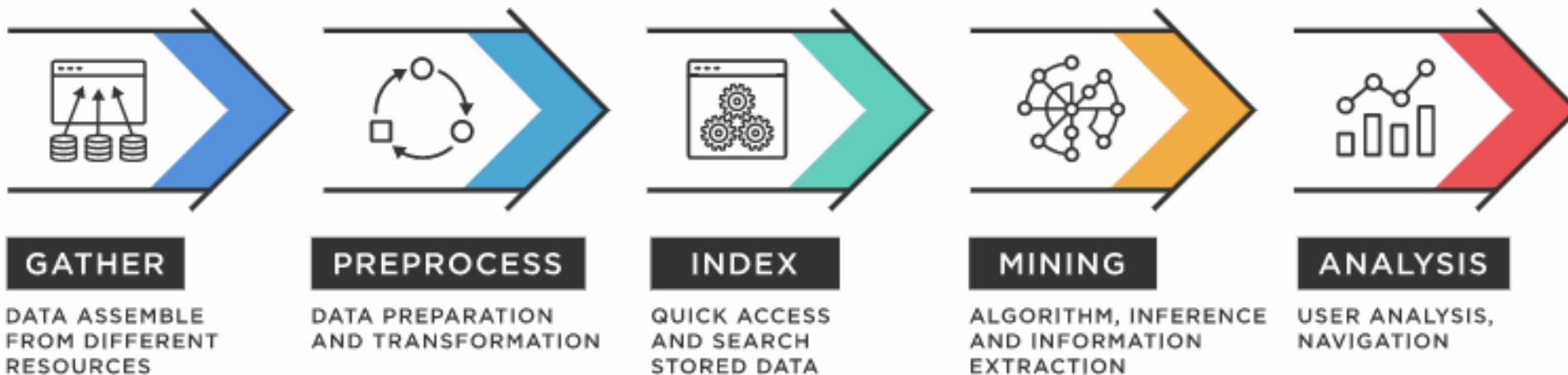
- Text: Data Science Lab is Yonsei University's best academic society.
- Text: 데이터 사이언스 랩 이즈 연세이 유니버시티즈 베스트 애커데믹 소우사이티. (Data Science Lab is Yonsei University's best academic society.)
- Buttons: 번역 수정 (Edit translation), 번역 평가 (Evaluate translation)

1. Text Mining



What is Text Mining?

텍스트 마이닝이란 어떤 글 속에서 **중요한 정보**를 뽑아내고 그 정보를 이용하여 **다양한 목적**으로 사용하는 것



1. Text Mining

YONSEI Data Science Lab | DSL

Crawling

BeautifulSoup

 Selenium

1. Text Mining

Cleansing & Tokenization

Tokenization : 주어진 코퍼스(corpus)에서 토큰(token)이라 불리는 단위로 나눈 작업

Cleansing : 갖고 있는 코퍼스(corpus)로부터 노이즈 데이터(오타, 특수문자, stop words 등)를 제거하는 것

Normalization : 표현 방법이 다른 단어들을 통합시켜서 같은 단어로 만들어 주는 것

1. Text Mining

Cleansing & Tokenization

ex) 데이터-사이언스-랩은 “텍스트 마이닝”에 대해서도 공부할 수 있는 좋은 학회입니다.

Tokenization : 주어진 코퍼스(corpus)에서 토큰(token)이라 불리는 단위로 나눈 작업

=> 형태소 tokenizer : 데이터 / 사이언스 / 랩 / 은 / 텍스트 / 마 / 이닝 / 에 / 대해 / 서도 /
공부 / 할 / 수 / 있는 / 좋은 / 학회 / 입니다 / .

=> 음절 단위 tokenizer : 데 / 이 / 터 / 사 / ... / 학 / 회 / 입 / 니 / 다 / .

Cleansing : 갖고 있는 코퍼스(corpus)로부터 노이즈 데이터(오타, 특수문자, stop words 등)를 제거하는 것

=> 데이터 사이언스 랩은 텍스트 마이닝에 대해서도 공부할 수 있는 좋은 학회입니다

Normalization : 표현 방법이 다른 단어들을 통합시켜서 같은 단어로 만들어 주는 것

=> 할, 하고, 하는, 함 - 하다 / 좋은, 좋다, 좋 - 좋다 / 입니다, 니다 - 이다 / ㅋㅋㅋㅋㅋㅋ, ㅋㅋㅋㅋㅋㅋㅋㅋㅋㅋ - ㅋㅋㅋ

Korean Tokenizer -konlpy

1. Kkma

['데이터', '사이언스', '랩', '은', '텍스트', '마이닝', '에', '대하', '어서', '도', '공부', '하', '르', '수', '있', '는', '중', '은', '학회', '이', '바니다', '.']

2. Komoran

['데이터', '사이언스', '랩', '은', '텍스트', '마', '이닝', '에', '대하', '아서', '도', '공부', '하', '르', '수', '있', '는', '중', '은', '학회', '이', '바니다', '.']

3. Hannanum

['데이터', '사이언스', '랩', '은', '텍스트', '마이닝', '에', '대', '어', '하', '어서', '도', '공부', '하', '르', '수', '있', '는', '중', '은', '학회', '이', '바니다', '.']

4. Okt

['데이터', '사이언스', '랩', '은', '텍스트', '마', '이닝', '에', '대해', '서도', '공부', '할', '수', '있는', '좋은', '학회', '입니다', '.']

5. Mecab

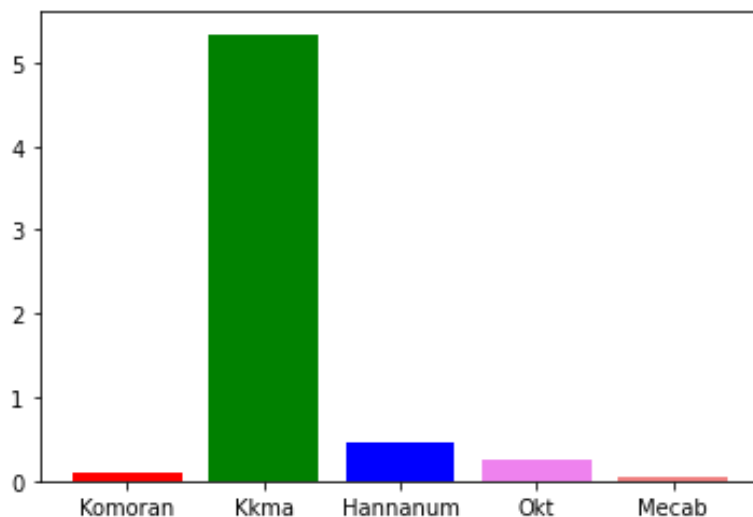
['데이터', '사이언스', '랩', '은', '텍스트', '마이닝', '에', '대해서', '도', '공부', '할', '수', '있', '는', '중', '은', '학회', '입니다', '.']

1. Text Mining

YONSEI Data Science Lab | DSL

Korean Tokenizer

형태소 분리 시 소요 시간



Okt 만의 특별한 기능

ex) 데사랩 진짜 좋은 것 같악ㅋㅋㅋㅋㅋㅋㅋ!

1. Stem : 단어의 어간을 추출해주는 기능
2. Norm : 문장을 정규화 해주는 기능

- Stem = False, Norm = False

('데', 'Noun'), ('싸', 'Verb'), ('랩', 'Noun'), ('진짜', 'Noun'), ('좋은', 'Adjective'), ('것', 'Noun'), ('같악', 'Noun'), ('ㅋㅋㅋㅋㅋㅋ', 'KoreanParticle'), ('!', 'Punctuation')

- Stem = True, Norm = False

('데', 'Noun'), ('싸다', 'Verb'), ('랩', 'Noun'), ('진짜', 'Noun'), ('좋다', 'Adjective'), ('것', 'Noun'), ('같악', 'Noun'), ('ㅋㅋㅋㅋㅋㅋ', 'KoreanParticle'), ('!', 'Punctuation')

- Stem = False, Norm = True

('데', 'Noun'), ('싸', 'Verb'), ('랩', 'Noun'), ('진짜', 'Noun'), ('좋은', 'Adjective'), ('것', 'Noun'), ('같아', 'Adjective'), ('ㅋㅋ', 'KoreanParticle'), ('!', 'Punctuation')

- Stem = True, Norm = True

('데', 'Noun'), ('싸다', 'Verb'), ('랩', 'Noun'), ('진짜', 'Noun'), ('좋다', 'Adjective'), ('것', 'Noun'), ('같다', 'Adjective'), ('ㅋㅋ', 'KoreanParticle'), ('!', 'Punctuation')

1. Text Mining

YONSEI Data Science Lab | DSL

Cleansing : 갖고 있는 코퍼스(corpus)로부터 노이즈 데이터(오타, 특수문자, stop words 등)를 제거하는 것

오타 제거

특수문자 제거

의미 없는 단어 (stop words) 제거 ex. 은, 는, 이, 가, 다, 이다, ㅋㅋㅋㅋ ...

! or ? : 특수문자네? 무조건 제거해야지!

ㅋㅋㅋㅋㅋㅋ : 제거해야 할까?

하고자 하는 task에 맞게끔..!

1. Text Mining

YONSEI Data Science Lab | DSL

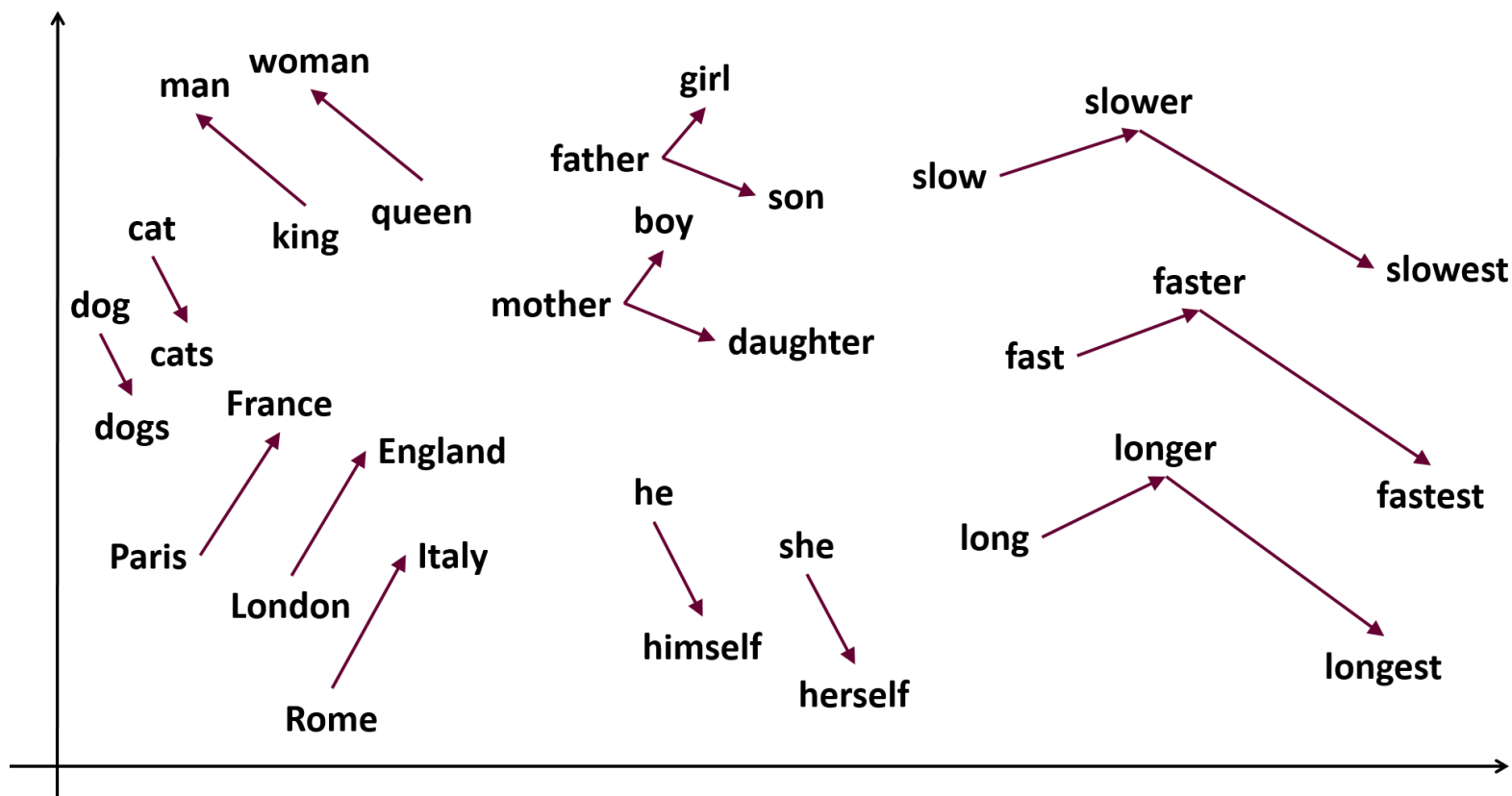
Tokenization & Cleansing

정답은 없습니다! **하고자 하는 task에 맞게끔..!**

다양한 시도와 실험 속에서 유의미한 정보를 찾는 것이 중요합니다!

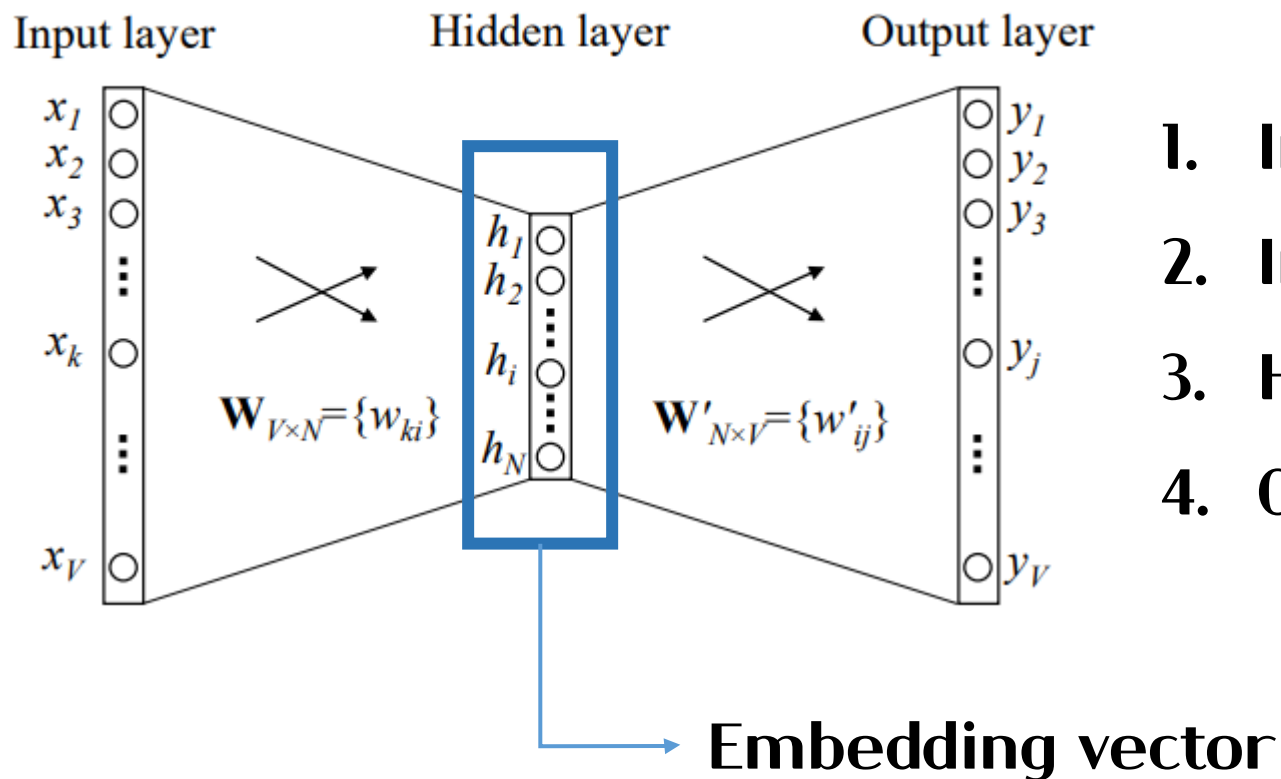
2. Word Embedding

word2vec : 단어 벡터 간 **유의미한 유사도**를 반영하도록 단어의 의미를 **수치화**하는 것



2. Word Embedding

word2vec : 단어 벡터 간 **유의미한 유사도**를 반영하도록 단어의 의미를 **수치화**하는 것



1. Input Data : One Hot Encoded Vector
2. Input \rightarrow Hidden : Embedding Matrix
3. Hidden \rightarrow Output : Score Matrix
4. Output \rightarrow Softmax function

2. Word Embedding

Corpus Tokenization & One Hot Encoding

Tokenization

우리집 강아지는 복슬 강아지 → 우리 / 집 / 강아지 / 는 / 복슬 / 강아지

One Hot Encoding

우리 : [1, 0, 0, 0, 0]

집 : [0, 1, 0, 0, 0]

강아지 : [0, 0, 1, 0, 0]

는 : [0, 0, 0, 1, 0]

복슬 : [0, 0, 0, 0, 1]

2. Word Embedding

word2vec - CBoW vs SkipGram

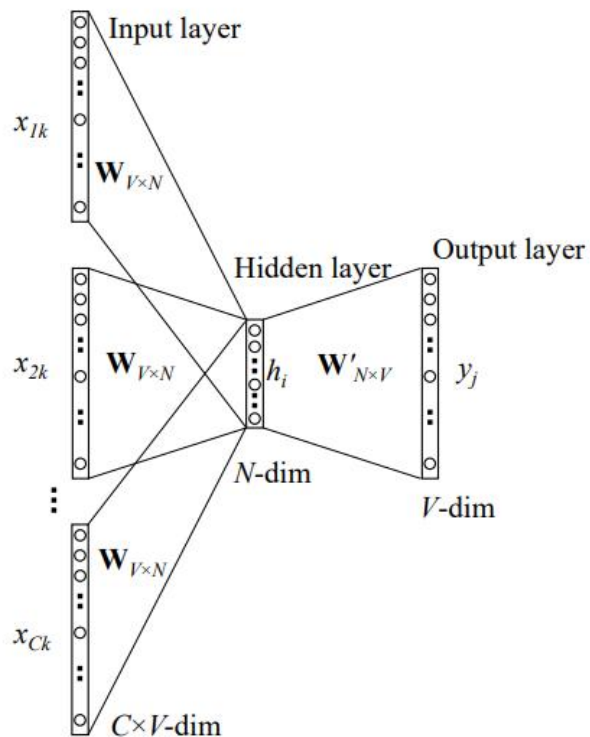


Figure 2: Continuous bag-of-words model

우리 집 ___는 복슬 강아지

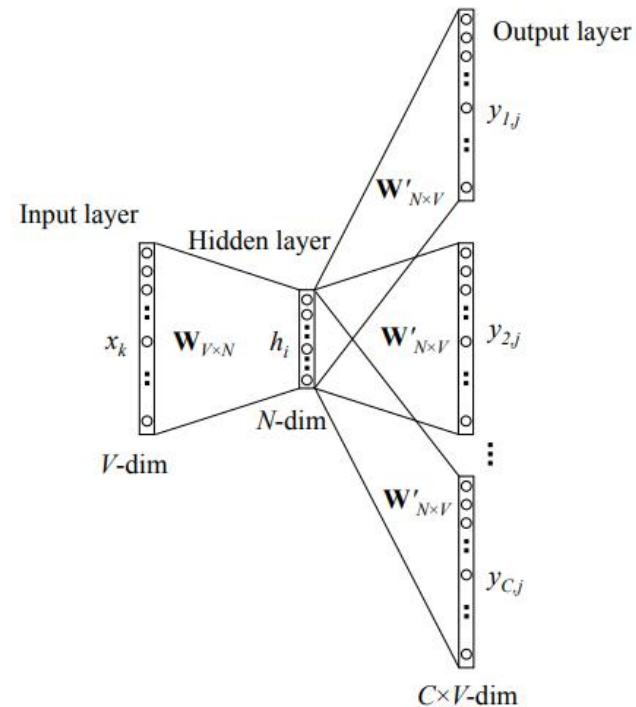


Figure 3: The skip-gram model.

___ 강아지_ _ _ _ _

2. Word Embedding

word2vec - CBoW vs SkipGram

CBoW

Input	Output
boy, is, going	the
the, is, going, to	boy
the, boy, going, to school	is
the, boy, is, to, school	going
boy, is, going, school	to
is, going, to	school

the	1
boy	1
is	1
going	1
to	1
school	1

Skip-Gram

Input	Output
the	boy, is, going
boy	the, is, going, to
is	the, boy, going, to school
going	the, boy, is, to, school
to	boy, is, going, school
school	is, going, to

the	3
boy	4
is	5
going	5
to	5
school	3

SkipGram이 CBoW보다 같은 epoch으로 학습을 한다면
각 단어들은 여러 번 여러 context에 걸쳐 빈번하게 학습됨.

성능 : SkipGram > CBoW

2. Word Embedding

word2vec - SkipGram 학습과정 살펴보기

Hyperparameter !

Hidden layer의 차원 (N) : 몇 차원의 벡터로 단어를 표현할 것인가

Window Size (m) : 주변 몇 개의 단어를 사용하여 context를 파악할 것인가

Goal !

Maximize $P(w_{c-m}, \dots, w_{c+m} | w_c)$

=> Loss function = $-\log P(w_{c-m}, \dots, w_{c+m} | w_c)$

Example

우리 집 강아지는 복슬 강아지 (_ _ _ 강아지 _ _ _ _ _)

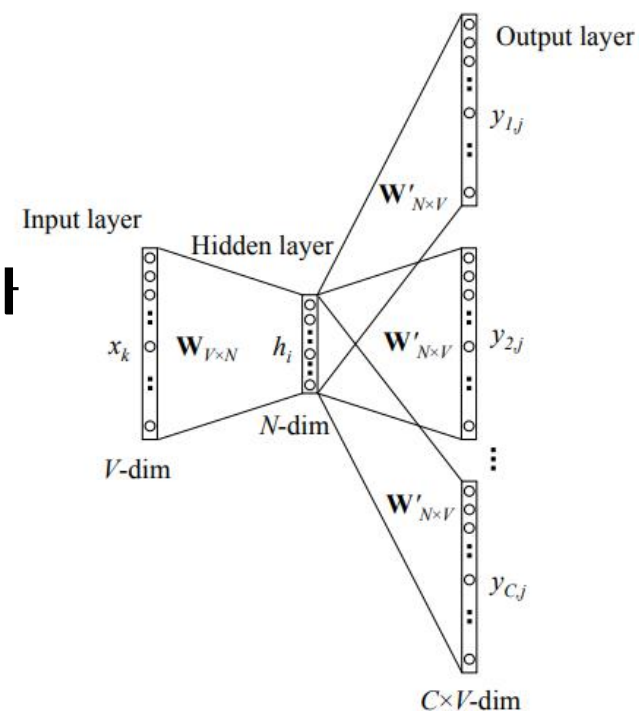


Figure 3: The skip-gram model.

2. Word Embedding

word2vec - SkipGram 학습과정 살펴보기 (Input Layer)

Tokenization

우리집 강아지는 복슬 강아지

=> 우리 / 집 / 강아지 / 는 / 복슬 / 강아지

One Hot Encoding

우리 : [1, 0, 0, 0, 0]

집 : [0, 1, 0, 0, 0]

강아지 : [0, 0, 1, 0, 0]

는 : [0, 0, 0, 1, 0]

복슬 : [0, 0, 0, 0, 1]

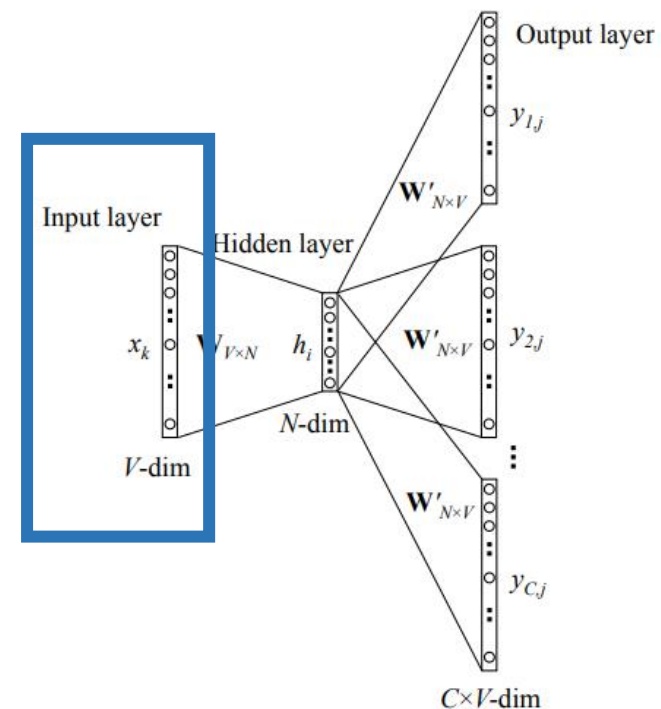


Figure 3: The skip-gram model.

Example

우리 집 강아지는 복슬 강아지 (__ _ 강아지 _ _ _ _)

2. Word Embedding

word2vec - SkipGram 학습과정 살펴보기 (Input Layer → Hidden Layer)

One Hot Encoding 한 vector와 Embedding matrix를 곱해,
N 차원의 embedding vector를 얻는다.

강아지 : [0, 0, 1, 0, 0] @ Embedding matrix

⇒ 강아지의 embedding vector : [0.2, 0.8, 1.9, -0.5]

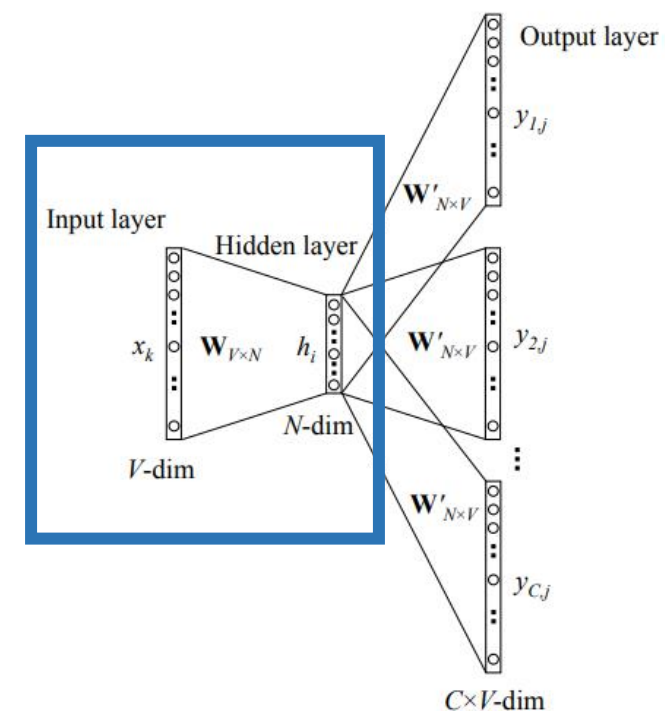


Figure 3: The skip-gram model.

Example

우리 집 강아지는 복슬 강아지 (__ _ 강아지 _ _ _ _)

2. Word Embedding

word2vec - SkipGram 학습과정 살펴보기 (Hidden Layer → Output Layer)

N 차원의 embedding vector와 Score matrix를 곱해
 $C(=2m)$ 개 만큼 output vector (score vector)를 얻는다.

강아지 : [0.2, 0.8, 1.9, -0.5] @ Score matrix

⇒ 강아지의 score vector : [0.5, 0.1, -0.8, -1.5, 0.9] × C

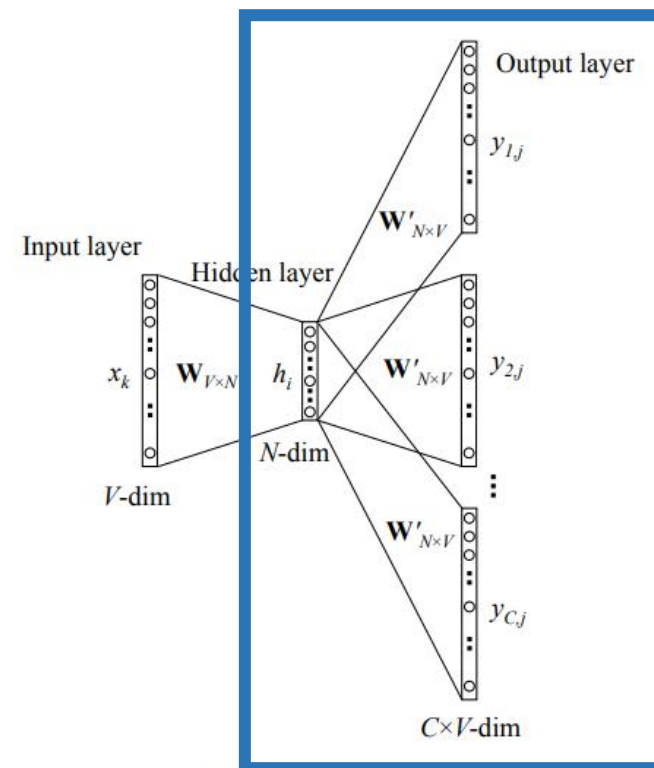


Figure 3: The skip-gram model.

Example

우리 집 강아지는 복슬 강아지 (__ _ 강아지 _ _ _)

2. Word Embedding

word2vec - SkipGram 학습과정 살펴보기 (Output Layer)

Score vector를 확률값으로 나타내기 위해 Softmax 함수를 취한다.

그 후 target vector와 차이를 이용하여 역전파 알고리즘 사용하여 Embedding Matrix, Score Matrix 를 update한다.

$[0.5, 0.1, -0.8, -1.5, 0.9] \Rightarrow \text{Predict} : [0.28, 0.19, 0.08, 0.04, 0.41]$

Target : 우리 : $[1, 0, 0, 0, 0]$

집 : $[0, 1, 0, 0, 0]$

는 : $[0, 0, 0, 1, 0]$

복슬 : $[0, 0, 0, 0, 1]$

강아지 : $[0, 0, 1, 0, 0]$

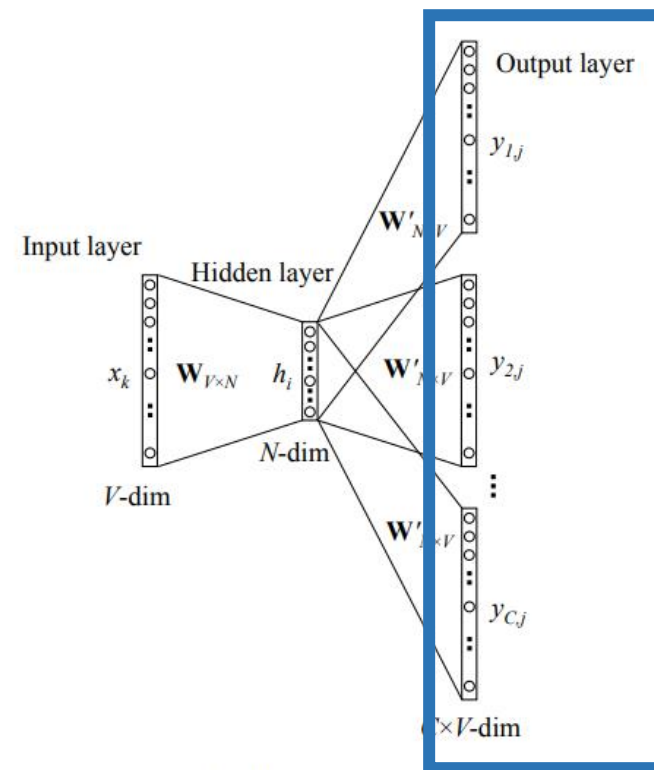


Figure 3: The skip-gram model.

Example

우리 집 강아지는 복슬 강아지 (__ _ 강아지 _ _ _)

2. Word Embedding

word2vec : 단어 벡터 간 **유의미한 유사도**를 반영하도록 단어의 의미를 **수치화**하는 것

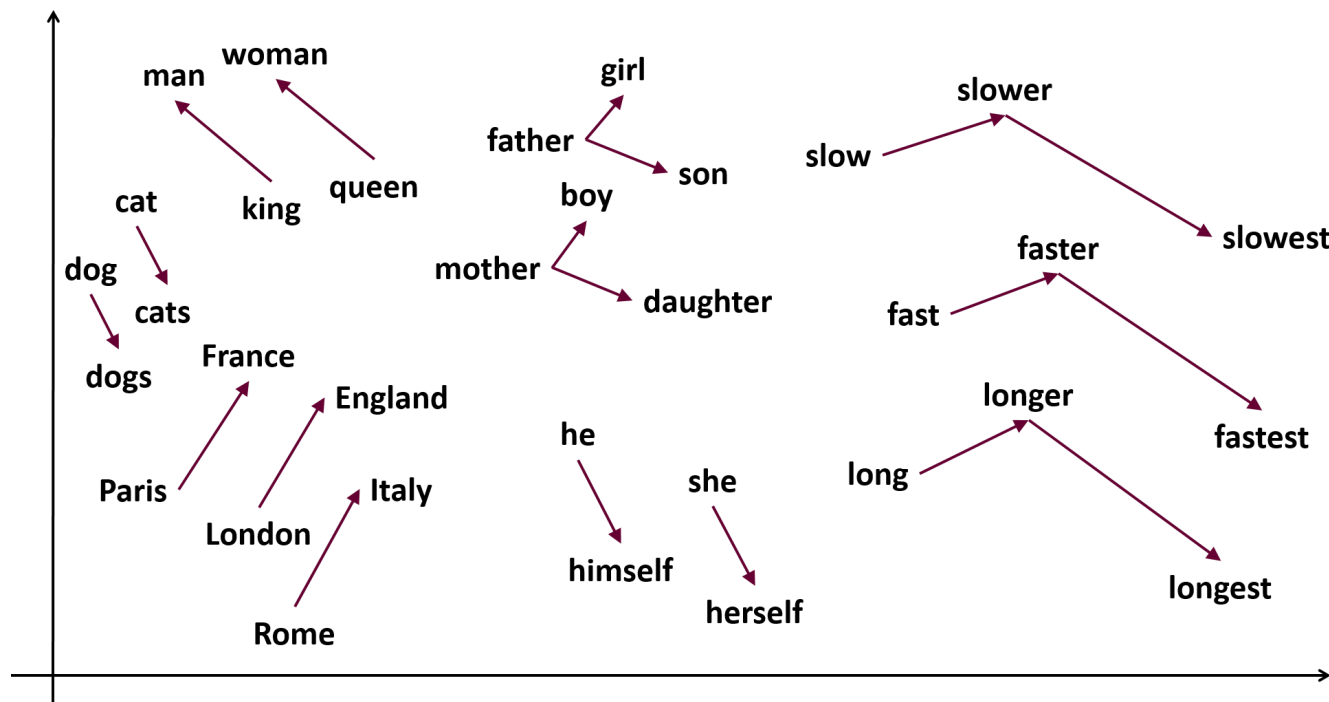
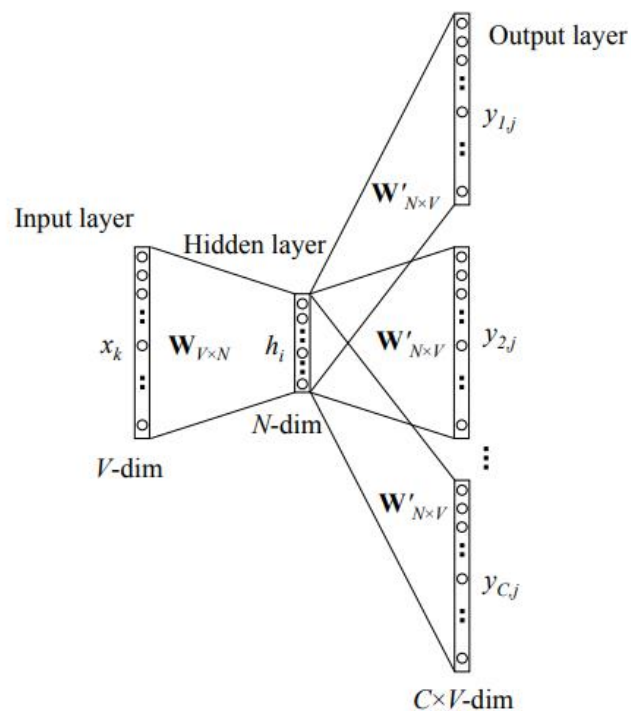
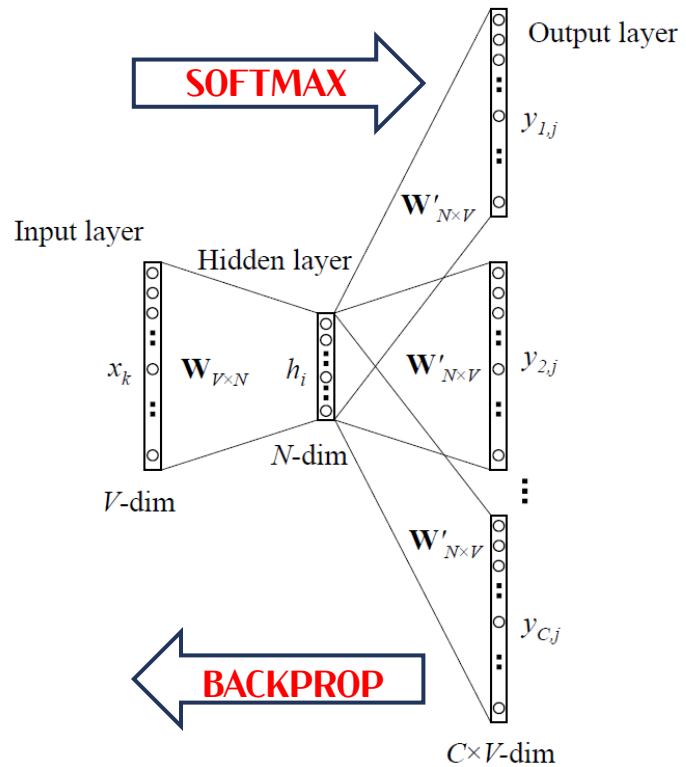


Figure 3: The skip-gram model.

2. Word Embedding

word2vec 단점 - Too Expensive Computation

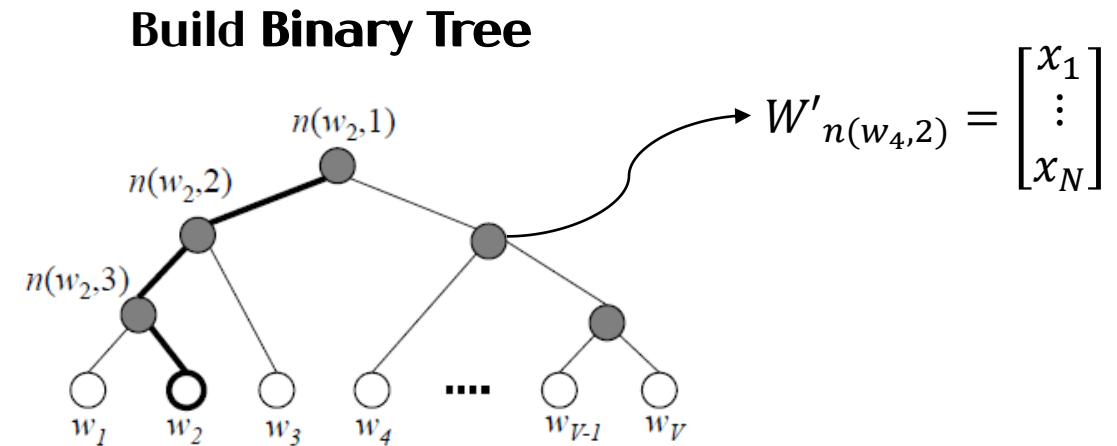
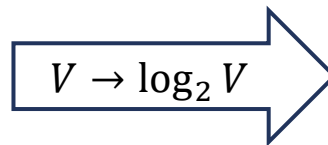
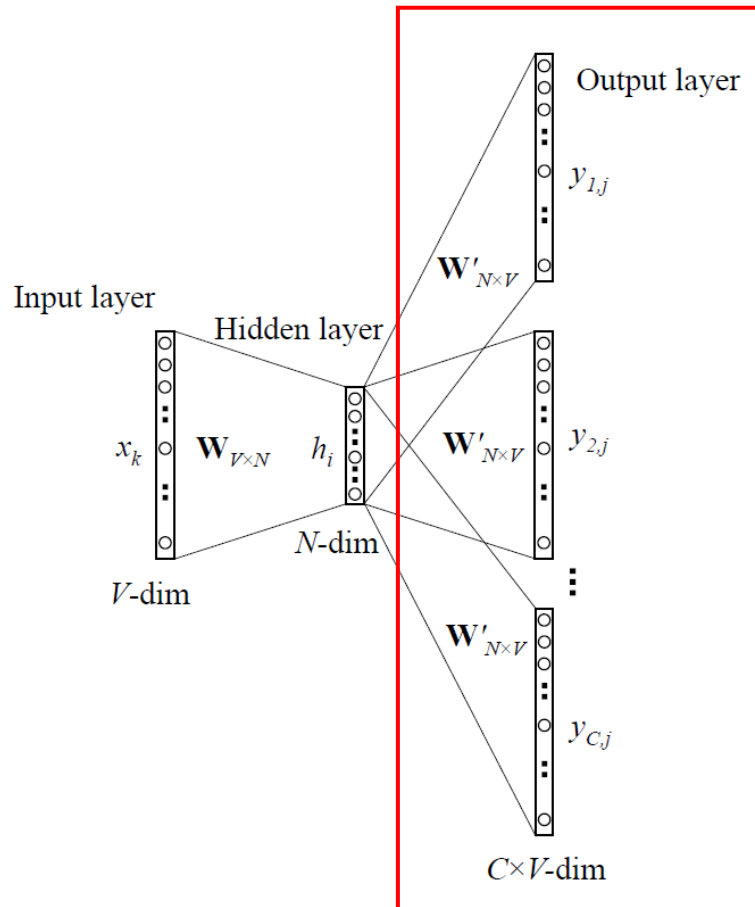


```
for each  $v_j \in W$  do (each word in context words)
  for each  $u_k \in \text{Context Words}$  do
     $EH \leftarrow 0$ 
    for each  $i = 1$  to  $V$  do
       $El_i \leftarrow \text{softmax}(W'^T v_{W_j}) - t_i$ 
       $EH \leftarrow EH + W' \cdot El_i$ 
       $v_{W'_i}^{new} \leftarrow v_{W'_i}^{old} - \alpha \cdot El_i \cdot h$ 
    end for
     $v_{W_j}^{new} \leftarrow v_{W_j}^{old} - \alpha \cdot EH^T$ 
  end for
end for
```

➔ Too expensive computation

2. Word Embedding

word2vec - Hierarchical Softmax

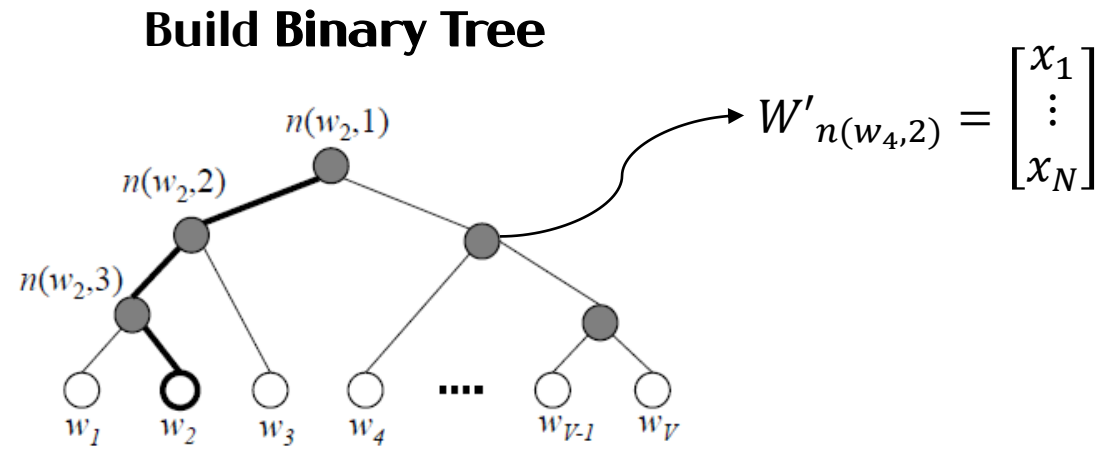


Build Binary Tree : V leaves, $V - 1$ inner units

$$p(w = w_o) = \prod_{j=1}^{L(w_o)-1} \sigma(\mathbb{I}[n(w_o, j+1) = ch(n(w_o, j))] \cdot W'^T_{n(w, j)} \cdot h)$$

2. Word Embedding

word2vec - Hierarchical Softmax



Build Binary Tree : V leaves, $V - 1$ inner units

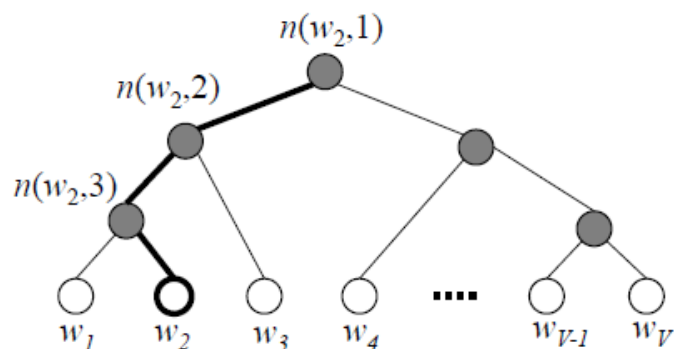
$$p(w = w_o) = \prod_{j=1}^{L(w_o)-1} \sigma(\mathbb{I}[n(w_o, j+1) = ch(n(w_o, j))] \cdot W'^T_{n(w, j)} \cdot h)$$

2. Word Embedding

word2vec - Hierarchical Softmax

$$p(w = w_o) = \prod_{j=1}^{L(w_o)-1} \sigma(\llbracket n(w_o, j+1) = ch(n(w_o, j)) \rrbracket) \cdot W_{n(w, j)}'^T \cdot h$$

Build Binary Tree



$\llbracket x \rrbracket = 1$ if x is True

$\llbracket x \rrbracket = -1$ if x is False

$ch(n)$: left child of unit n

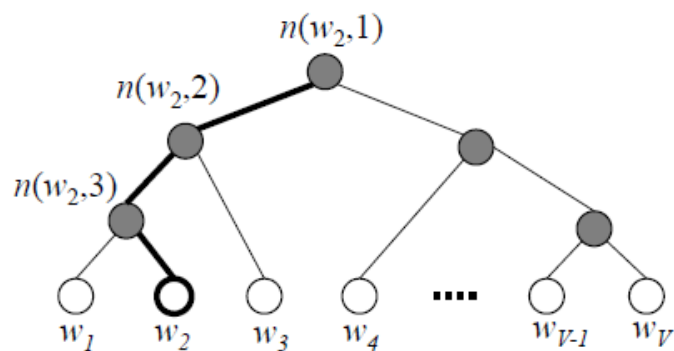
$n(w_o, j+1) = ch(n(w_o, j))$: w_o 까지 가는 경로에서
왼쪽으로 뻗어 나가면 True, 오른쪽으로 뻗어 나가면 False

Ex) w_2 까지 가는 경로에서 1, 1, -1 을 반환하게 됨.

2. Word Embedding

word2vec - Hierarchical Softmax

Build Binary Tree



$$p(w = w_o) = \prod_{j=1}^{L(w_o)-1} \sigma(\mathbb{I}[n(w_o, j+1) = ch(n(w_o, j))] \cdot W'_{n(w, j)}^T \cdot h)$$

$$p(w = w_2) = p(n(w_2, 1), left) \cdot p(n(w_2, 2), left) \cdot p(n(w_2, 3), right)$$

$$= \sigma(W'_{n(w_2, 1)}^T \cdot h) \cdot \sigma(W'_{n(w_2, 2)}^T \cdot h) \cdot \sigma(-W'_{n(w_2, 3)}^T \cdot h)$$

=> 연산량 : 3

$$p(w = w_2) = \frac{\exp(w_2)}{\sum_{j=1}^V \exp(w_j)}$$

=> 연산량 : V

2. Word Embedding

word2vec - Hierarchical Softmax

Error Function

Softmax

$$E = -\log \prod_{c=1}^C \frac{\exp(u_{c,j_c^*})}{\sum_{j'=1}^V \exp(u_{j'})} \quad (u_j = W_j'^T \cdot h)$$

Hierarchical Softmax

$$E = -\log \prod_{j=1}^{L(w_o)-1} \sigma(\mathbb{I}(n(w_o, j+1) = ch(n(w_o, j))) \cdot W_{n(w,j)}'^T \cdot h)$$

Update W Matrix

Softmax

$$W_{w_I}^{new} = W_{w_I}^{old} - \alpha \cdot EH^T, \quad (EH_i = \sum_{j=1}^V E I_j W'_{ij})$$

Hierarchical Softmax

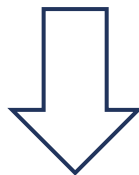
$$W_{w_I}^{new} = W_{w_I}^{old} - \alpha \cdot EH^T, \quad (EH = \sum_{j=1}^{L(w)-1} (\sigma(W_{n(w,j)}'^T \cdot h) - t_j) \cdot W'_{n(w,j)})$$

2. Word Embedding

word2vec – Negative Sampling

Only the weights corresponding to the target word might get a significant update.

The calculation of the final probabilities using the softmax is quite an expensive operation.



The idea of negative sampling is more straightforward than hierarchical softmax: in order to deal with the difficulty of having too many output vectors that need to be updated per iteration, we only update a sample of them.

Apparently the output word (i.e., the ground truth, or positive sample) should be kept in our sample and gets updated, and we need to sample a few words as negative samples (hence “negative sampling”). A probabilistic distribution is needed for the sampling process, and it can be arbitrarily chosen. We call this distribution the noise distribution, and denote it as $P_n(w)$. One can determine a good distribution empirically.⁶

2. Word Embedding

word2vec - Negative Sampling

The idea of negative sampling is more straightforward than hierarchical softmax: in order to deal with the difficulty of having too many output vectors that need to be updated per iteration, we only update a sample of them.

Apparently the output word (i.e., the ground truth, or positive sample) should be kept in our sample and gets updated, and we need to sample a few words as negative samples (hence “negative sampling”). A probabilistic distribution is needed for the sampling process, and it can be arbitrarily chosen. We call this distribution the noise distribution, and denote it as $P_n(w)$. One can determine a good distribution empirically.⁶

$$E = -\log \sigma(\mathbf{v}'_{w_O}{}^T \mathbf{h}) - \sum_{w_j \in \mathcal{W}_{\text{neg}}} \log \sigma(-\mathbf{v}'_{w_j}{}^T \mathbf{h})$$

$$W_{\text{neg}} = \{w_j | j = 1, 2, 3, \dots, k\} \sim P_n(w): \text{noise dist'n}$$

$$P_n(w_i) = \frac{f(w_i)}{\sum_{j=0}^n f(w_j)}, \quad f(w_i): \text{frequency of } w_i, \quad \text{for better result, } P_n(w_i) = \frac{f(w_i)^{\frac{3}{4}}}{\sum_{j=0}^n f(w_j)^{\frac{3}{4}}}$$

$$0.9^{\frac{3}{4}} = 0.9240, \quad 0.2^{\frac{3}{4}} = 0.2990$$

2. Word Embedding

word2vec - Negative Sampling

$$\begin{aligned}\theta &= \operatorname{argmax}_{\theta} \prod_{(w,c) \in D} P(D=1|w,c,\theta) \prod_{(w,c) \in \tilde{D}} P(D=0|w,c,\theta) \\ &= \operatorname{argmax}_{\theta} \prod_{(w,c) \in D} P(D=1|w,c,\theta) \prod_{(w,c) \in \tilde{D}} (1 - P(D=1|w,c,\theta)) \\ &= \operatorname{argmax}_{\theta} \sum_{(w,c) \in D} \log P(D=1|w,c,\theta) + \sum_{(w,c) \in \tilde{D}} \log(1 - P(D=1|w,c,\theta)) \\ &= \operatorname{argmax}_{\theta} \sum_{(w,c) \in D} \log \frac{1}{1 + \exp(-u_w^T v_c)} + \sum_{(w,c) \in \tilde{D}} \log \left(1 - \frac{1}{1 + \exp(-u_w^T v_c)}\right) \\ &= \operatorname{argmax}_{\theta} \underbrace{\sum_{(w,c) \in D} \log \frac{1}{1 + \exp(-u_w^T v_c)}}_{\text{Maximize}} + \underbrace{\sum_{(w,c) \in \tilde{D}} \log \left(\frac{1}{1 + \exp(u_w^T v_c)}\right)}_{\text{Minimize}}\end{aligned}$$

Maximize the probability of co-occurrence for actual words that lie in the context

Minimize the probability of co-occurrence for some random words that don't lie in the context

$$O = \underbrace{\log \sigma(\mathbf{v}'_{w_O}{}^T \mathbf{h})}_{\text{Maximize}} + \sum_{w_j \in \mathcal{W}_{\text{neg}}} \log \sigma(-\mathbf{v}'_{w_j}{}^T \mathbf{h}) \quad E = - \underbrace{\log \sigma(\mathbf{v}'_{w_O}{}^T \mathbf{h})}_{\text{Minimize}} - \sum_{w_j \in \mathcal{W}_{\text{neg}}} \log \sigma(-\mathbf{v}'_{w_j}{}^T \mathbf{h})$$

2. Word Embedding

word2vec - Negative Sampling

Ex. 우리 집 강아지는 복슬 강아지



Input	Target
강아지	집
강아지	는
는	강아지
는	복슬
복슬	는
복슬	강아지

Input1	Input2	Output
강아지	집	1
강아지	는	1
는	강아지	1
는	복슬	1
복슬	는	1
복슬	강아지	1

2. Word Embedding

word2vec - Negative Sampling

Ex. 우리 집 강아지는 복슬 강아지



Input1	Input2	Output
강아지	집	1
강아지	는	1
는	강아지	1
는	복슬	1
복슬	는	1
복슬	강아지	1

Input1	Input2	Output
강아지	집	1
강아지	는	1
강아지	우리	0
는	강아지	1
는	복슬	1
는	우리	0
복슬	는	1
복슬	강아지	1
복슬	집	0

Randomly choose from $P_n(w)$ →

Randomly choose from $P_n(w)$ →

Randomly choose from $P_n(w)$ →

2. Word Embedding

word2vec - Negative Sampling

Ex. 우리 집 강아지는 복슬 강아지



$$\text{minimize } E = -\log \sigma(\mathbf{v}'_{w_o}{}^T \mathbf{h}) - \sum_{w_j \in \mathcal{W}_{\text{neg}}} \log \sigma(-\mathbf{v}'_{w_j}{}^T \mathbf{h})$$

Input1	Input2	Output
강아지	집	1
강아지	는	1
강아지	우리	0
는	강아지	1
는	복슬	1
는	우리	0
복슬	는	1
복슬	강아지	1
복슬	집	0

2. Word Embedding

word2vec - Negative Sampling

Error Function

Softmax

$$E = -\log \prod_{c=1}^C \frac{\exp(u_{c,j_c^*})}{\sum_{j'=1}^V \exp(u_{j'})} \quad (u_j = v'_{w_j} \cdot h)$$

Negative Sampling

$$E = -\log \sigma(v'_{w_o} \cdot h) - \sum_{w_j \in W_{neg}} \log \sigma(-v'_{w_j} \cdot h)$$

Backpropagation

Softmax

$$v_{w_I}^{new} = v_{w_I}^{old} - \alpha \cdot EH^T, \quad (EH_i = \sum_{j=1}^V EI_j v'_{ij})$$

Negative Sampling

$$v_{w_I}^{new} = v_{w_I}^{old} - \alpha \cdot EH^T, \quad (EH = \sum_{w_j \in \{w_o\} \cup W_{neg}} (\sigma(v'_{w_j} \cdot h) - t_j) \cdot v'_{w_j})$$

2. Word Embedding

Other Methods

FastText : word2vec의 확장판.

1. 내부단어(subword)의 학습

자연어처리 -> <자연, 자연어, 언어처, 어처리, 처리>
-> <자, 자_, _어, _어_, _어_, ...>

2. 모르는 단어(Out Of Vocabulary, OOV)에 대한 대응

3. 단어 집합 내 빈도 수가 적었던 단어(Rare Word)에 대한 반응

GloVe : word2vec의 예측 기반 방법론과 카운트 기반 방법론을 모두 사용함.

3. Visualization

PCA

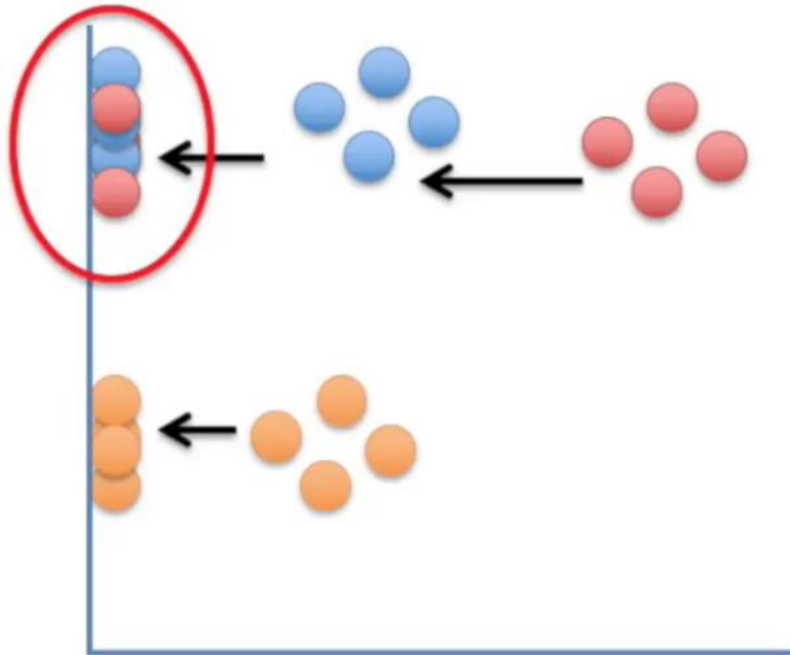
Dimension Reduction : N차원으로 나타낸 embedding vector를 눈으로 보고 싶다!

설명력이 높은 2~3개의 주성분을 이용하여 차원 축소 후 plot!

옛날에 배웠으니까 넘어갈게요~!

3. Visualization

PCA



PCA의 단점

선형 분석 방식으로 값을 mapping 하기 때문에
차원이 감소하면서 군집화 되어 있는 데이터들이 뭉개져서
제대로 구별할 수 없는 문제가 발생할 수 있음.

3. Visualization

t-SNE (t-distributed Stochastic Neighbor Embedding)

$$p_{j|i} = \frac{e^{-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma_i^2}}}{\sum_{k \neq i} e^{-\frac{\|\mathbf{x}_i - \mathbf{x}_k\|^2}{2\sigma_i^2}}}$$

원래 차원에서 i가 j를 이웃으로 선택할 확률

$$q_{j|i} = \frac{e^{-\|\mathbf{y}_i - \mathbf{y}_j\|^2}}{\sum_{k \neq i} e^{-\|\mathbf{y}_i - \mathbf{y}_k\|^2}}$$

축소한 차원에서 i가 j를 이웃으로 선택할 확률

y를 구할 수 있는가? 지금 당장은 못 구한다...

$$\Rightarrow Cost = \sum_i KL(P_i || Q_i) = \sum_i \sum_j p_{j|i} \log \frac{p_{j|i}}{q_{j|i}}$$

$$\frac{\partial C}{\partial \mathbf{y}_i} = 2 \sum_j (\mathbf{y}_j - \mathbf{y}_i) (p_{j|i} - q_{j|i} + p_{i|j} - q_{i|j}) \quad \text{!!!Update rule!!!}$$

3. Visualization

t-SNE (t-distributed Stochastic Neighbor Embedding)

$$p_{j|i} = \frac{e^{-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma_i^2}}}{\sum_{k \neq i} e^{-\frac{\|\mathbf{x}_i - \mathbf{x}_k\|^2}{2\sigma_i^2}}}$$

원래 차원에서 i가 j를 이웃으로 선택할 확률

$$q_{j|i} = \frac{e^{-\|\mathbf{y}_i - \mathbf{y}_j\|^2}}{\sum_{k \neq i} e^{-\|\mathbf{y}_i - \mathbf{y}_k\|^2}}$$

축소한 차원에서 i가 j를 이웃으로 선택할 확률

$$p_{ij} = \frac{e^{-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma_i^2}}}{\sum_{k \neq l} e^{-\frac{\|\mathbf{x}_k - \mathbf{x}_l\|^2}{2\sigma_i^2}}} \rightarrow p_{ij} = \frac{p_{j|i} + p_{i|j}}{2n} \quad \sum_j p_{ij} > \frac{1}{2n}$$

i와 j가 서로를 이웃으로 선택할 확률 (pairwise probability)

$$Cost = \sum_i KL(P_i || Q_i) = \sum_i \sum_j p_{j|i} \log \frac{p_{j|i}}{q_{j|i}}$$

$$\frac{\partial C}{\partial \mathbf{y}_i} = 2 \sum_j (\mathbf{y}_j - \mathbf{y}_i) (p_{j|i} - q_{j|i} + p_{i|j} - q_{i|j})$$

$$Cost = \sum_i KL(P_i || Q_i) = \sum_i \sum_j p_{ij} \log \frac{p_{ij}}{q_{ij}}$$

$$\frac{\partial C}{\partial \mathbf{y}_i} = 4 \sum_j (\mathbf{y}_j - \mathbf{y}_i) (p_{ij} - q_{ij})$$

Perplexity : SNE는 거리에 반비례하게 영향력을 정의하여 embedding 하는데 사용한다.

이때 Perplexity란 **어느 범위까지 영향력을 강하게 할 것인가**를 결정하는 hyperparameter

5~50 사이의 값에서 robust함.

3. Visualization

t-SNE (t-distributed Stochastic Neighbor Embedding)

왜 “t”-SNE일까??

Crowding problem : $q_{j|i} = \frac{e^{-\|\mathbf{y}_i - \mathbf{y}_j\|^2}}{\sum_{k \neq i} e^{-\|\mathbf{y}_i - \mathbf{y}_k\|^2}}$

무슨 분포가 생각나시나요?

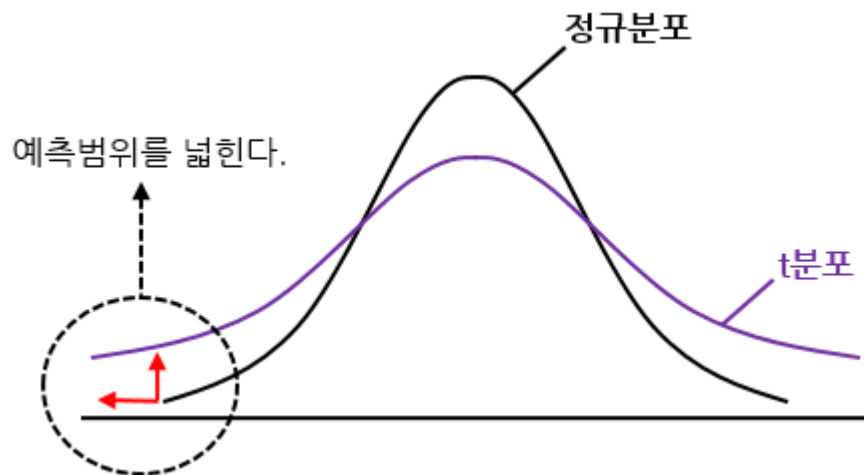
3. Visualization

t-SNE (t-distributed Stochastic Neighbor Embedding)

왜 “t”-SNE일까??

Crowding problem : 정규분포는 중심에서부터 멀어지면 급격히 값이 감소함.

어느 정도 멀리 있는 데이터들은 서로 간의 위치 정보를 담을 수 없음.

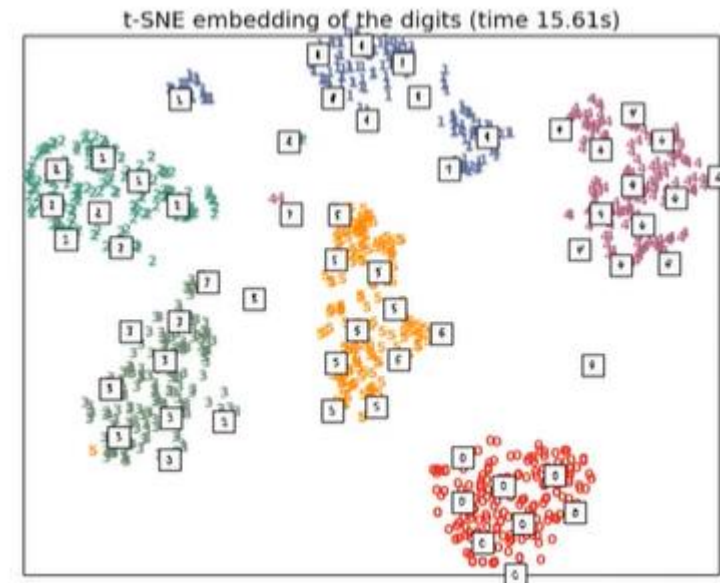
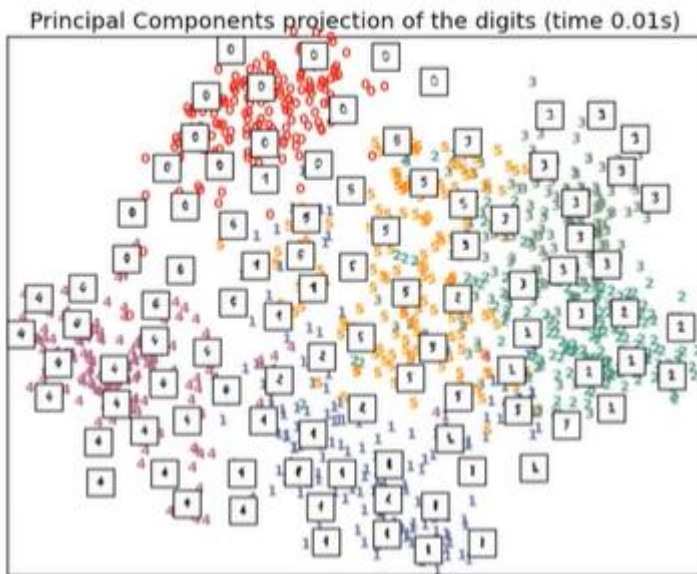


$$q_{j|i} = \frac{e^{-\|\mathbf{y}_i - \mathbf{y}_j\|^2}}{\sum_{k \neq i} e^{-\|\mathbf{y}_i - \mathbf{y}_k\|^2}} \Rightarrow q_{ji} = \frac{(1 + \|\mathbf{y}_i - \mathbf{y}_j\|^2)^{-1}}{\sum_{k \neq l} (1 + \|\mathbf{y}_k - \mathbf{y}_l\|^2)^{-1}}$$

$$\frac{\partial C}{\partial \mathbf{y}_i} = 4 \sum_j (\mathbf{y}_j - \mathbf{y}_i) (p_{ij} - q_{ij}) (1 + \|\mathbf{y}_i - \mathbf{y}_j\|^2)^{-1}$$

3. Visualization

t-SNE (t-distributed Stochastic Neighbor Embedding)



t-SNE를 사용 했을 때 PCA보다 더 군집이 잘 떨어져 있는 모습을 확인할 수 있음.

<https://distill.pub/2016/misread-tsne/>

1. Tokenization & Cleansing

1. 한국어의 경우 형태소 단위, 음절 단위 등으로 분석 가능함
2. KoNLPy + Huggingface 의 Tokenizer 사용 가능함
3. Tokenization이든 Cleansing에든 정해진 방법은 없음

2. Word Embedding - word2vec

1. word2vec은 좋은 word embedding 방법임
2. Heavy computation을 방지하기 위해 HS, NS가 고안됨
3. 그 외에도 FastText, GloVe 같은 embedding 방법이 있음

3. Visualization - PCA, t-SNE

1. PCA, t-SNE 모두 좋은 저차원 embedding 방법임
2. 여러 가지 시도해보고 좋은 것으로 visualization 하는 것이 좋음

Reference

https://en.wikipedia.org/wiki/Text_mining

딥 러닝을 이용한 자연어 처리 입문 <https://wikidocs.net/book/2155>

Rong, Xin. "word2vec parameter learning explained." (2014).

Mikolov, Tomas, et al. "Efficient estimation of word representations in vector space." (2013).

<https://www.youtube.com/watch?v=INHwh8k4XhM>

Image Reference

<https://www.tibco.com/ko/reference-center/what-is-text-mining>

<https://www.expressanalytics.com/blog/social-media-sentiment-analysis/>

<https://www.samyzaf.com/ML/nlp/nlp.html>



Thank you!!!

22.08.30 / DSL 7기 최명헌