# 자연어처리팀

2팀

한준호

윤세인

김나현

윤여원

권능주

# **INDEX**

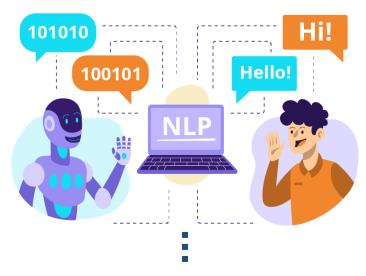
- 1. 자연어
- 2. 단어표현방법
- 3. 언어모델링
- 4. 기계번역

# 1

# 자연어

자연어

의식적인 계획이나 계획 없이 사용, 반복, 변화의 과정을 거쳐 인간 공동체에서 발생한 모든 언어



자연어 처리의 목표는 컴퓨터가 문서 내 언어의 문맥적 뉘앙스를 포함해 문서 내용을 이해할 수 있도록 하는 것

#### 자연어의 특징

#### 순차성

자연어의 가장 대표적인 특징으로 순서가 달라지는 경우 의미가 손상됨

#### 문맥의 중요성

언어가 가진 다의성으로 하나의 단어나 문장이 상황에 따라 다양한 의미로 활용

#### 문장 내 정보 부족

언어는 효율성을 극대화 하기 위해 많은 정보를 생략하는 경우가 빈번하게 발생

#### 자연어의 특징

순차성

자연어의 가장 대표적인 특징으로 순서가 달라지는 경우 의미가 손상됨

문맥의 중요성

나는 밥을 먹었다.

밥이 나를 먹었다.

호디 이프로 이미에도 이파기 보다에게 된 이

단어의 순서가 바뀜으로 문장의 의미가 달라짐

#### 자연어의 특징

문맥의 중요성

언어가 가진 다의성으로

하나의 단어나 문장이 상황에 따라 다양한 의미로 활용

순차성

배가 달달하니 맛있네!

과일의 한 종류 배 ㅣ 교통수단의 일종인 배

→ 뉘앙스를 통해 뜻을 파악하기 위해 문맥 해석이 중요

#### 자연어의 특징

#### 문장 내 정보 부족

언어는 효율성을 극대화 하기 위해 많은 정보를 생략하는 경우가 빈번하게 발생

문맥의 중요성

내가 한거 아니야!

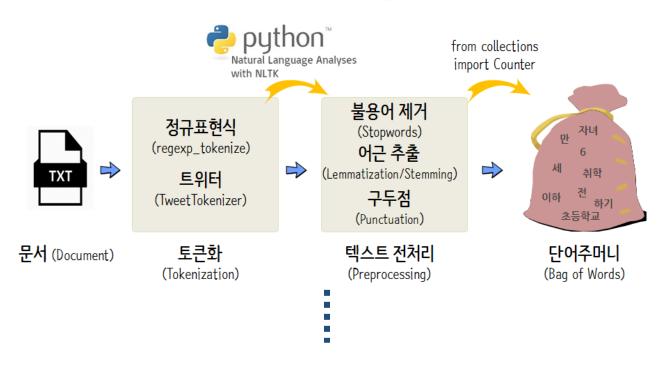
무엇을 하지 않았는지 생략

군시기 필니시는 경우 의미기 **는**경험

정보의 생략이 많을수록 자연어 처리가 어려워짐

#### 자연어 전처리

#### 자연어 처리의 절차



텍스트 데이터를 모델의 입력으로 사용하기 위해 자연어 특성을 반영할 수 있는 변환방법이 선행

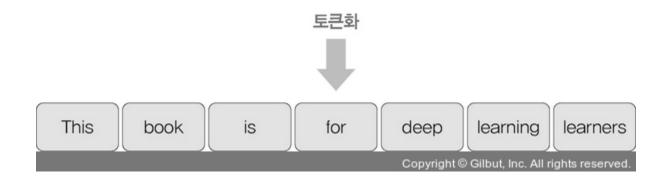
자연어 전처리 : 토큰화

토큰화

주어진 말뭉치(corpus)에서 토큰이라 불리는 단위로 나누는 작업

의미 있는 단위로 토큰을 정의

"This book is for deep learning learners"



자연어 전처리 : 토큰화

토큰화

주어진 말뭉치(corpus)에서 토<mark>큰</mark>이라 불리는 단위로 나누는 작업

의미 있는 단위로 토큰을 정의

"This book is for deep learning learners"

가장 기본적인 토큰화는 단어 토큰화

단어 외에도 단어구, 문장 토큰화도 가능

명확한 선택 기준이 없어 필요에 따라 적절한 토큰화 선택

자연어 전처리 : 토큰화 방법

"Hey Amazon - my package never arrived https://www.amazon.com/gp/css/order-history?ref\_=nav\_orders\_first PLEASE FIX ASAP! @amazonhelp"

세가지 토큰화 방법 사용



∰공백 기준으로 분리

['Hey', 'Amazon', '-', 'my', 'package', 'never', 'arrived', 'https://www.amazon.com/gp/css/order-history?ref\_=nav\_orders\_first', 'PLEASE', 'FIX', 'ASAP!', '@amazonhelp']

── 단순히 공백을 기준으로 분리

자연어 전처리 : 토큰화 방법

"Hey Amazon - my package never arrived https://www.amazon.com/gp/css/order-history?ref\_=nav\_orders\_first PLEASE FIX ASAP! @amazonhelp"

세가지 토큰화 방법 사용



≝keras의 text\_to\_word\_sequence 함수

['hey', 'amazon', 'my', 'package', 'never', 'arrived', 'https', 'www',
'amazon', 'com', 'gp', 'css', 'order', 'history', 'ref', 'nav', 'orders',
'first', 'please', 'fix', 'asap', 'amazonhelp']

──> 문장 부호와 특수문자를 제거하고 대문자를 소문자로 변환

자연어 전처리 : 토큰화 방법

"Hey Amazon - my package never arrived https://www.amazon.com/gp/css/order-history?ref\_=nav\_orders\_first PLEASE FIX ASAP! @amazonhelp"

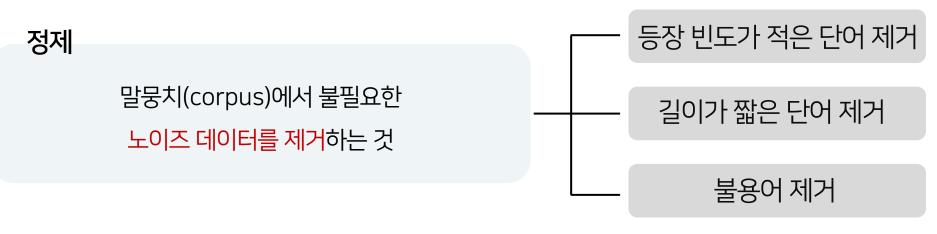
세가지 토큰화 방법 사용



∰nltk의 word\_tokenize 함수

['Hey', 'Amazon', '-', 'my', 'package', 'never', 'arrived', 'https', ':',
'//www.amazon.com/gp/css/order-history', '?', 'ref\_=nav\_orders\_first',
'PLEASE', 'FIX', 'ASAP', '!', '@', 'amazonhelp']

자연어 전처리 : 정제





토큰화 작업 전, 후에 텍스트 데이터를 용도에 맞게 정제 및 <mark>정규화</mark> 하는 작업을 진행

자연어 전처리: 정제

"Hey Amazon - my package never arrived https://www.amazon.com/gp/css/order-history?ref\_=nav\_orders\_first PLEASE FIX ASAP! @amazonhelp"



#### <del>일</del>등장 빈도가 적은 단어 제거

['Hey', 'Amazon', '-', 'my', 'package', 'never', 'arrived', 'PLEASE', 'FIX',
'ASAP', '!', '@', 'amazonhelp']

자연어 전처리: 정제

"Hey Amazon - my package never arrived https://www.amazon.com/gp/css/order-history?ref\_=nav\_orders\_first PLEASE FIX ASAP! @amazonhelp"



### ∰길이가 짧은 단어 제거

['Hey', 'Amazon', 'my', 'package', 'never', 'arrived', 'PLEASE', 'FIX', 'ASAP',
'amazonhelp']

**---->** 길이가 짧은 문장 부호와 특수 문자를 제거

자연어 전처리: 정제



#### 불용어

"Hey Amazon" - my package never arrived https://www.amazon.com/gp/css/order-historyPref\_=nav\_Arders first please HistoryPref\_=nav\_Arders first please HistoryP



['Hey', 'Amazon', 패키지에서 제공하는 <mark>불용어 목록</mark>에 추가로, 'PLEASE', 'FIX', 'ASAP', '!', '@ 중요하지 않다고 판단되는 단어를 지정해서 사용

자연어 전처리 : 정제

"Hey Amazon - my package never arrived https://www.amazon.com/gp/css/order-history?ref\_=nav\_orders\_first PLEASE FIX ASAP! @amazonhelp"



#### **뿔불용어 제거**

['hey', 'amazon', 'package', 'never', 'arrived', 'please', 'fix', 'asap',
'amazonhelp']

nltk 패키지의 stopwords 목록 사용

자연어 전처리 : 정규화

정규화

표현 방법이 다른 단어들을 통합해서 같은 단어로 만들어주는 것

### <mark>딸</mark>대소문자 통일

['hey', 'amazon', 'package', 'never', 'arrived', 'please', 'fix', 'asap',
'amazonhelp']

──> 대문자와 소문자는 의미 차이를 만들지 않는 경우 통일해서 사용

해당 예시의 경우 Amazon은 회사 명칭을 의미하지만 amazon은 다른 의미를 갖기 때문에 주의

자연어 전처리 : 정규화

#### ∰표제어 추출

['hey', 'amazon', 'package', 'never', 'arrive', 'please', 'fix', 'asap',
'amazonhelp']

nltk 패키지의 WordNetLemmatizer 사용

표제어 추출을 통해 arrived -> arrive로 변환

### ∰어간 추출

['hey', 'amazon', 'packag', 'never', 'arriv', 'pleas', 'fix', 'asap',
'amazonhelp']

nltk 패키지의 PorterStemmer 사용

자연어 전처리: 정규화



# 정제와 정규화

```
['hey', amazon', 토큰화:-> 정제'-> 정규화 과정을 거쳐se', 'fix', \asap',
'amazonhelp']
                  쓸만한 토큰을 생성하고
               필요에 따라 특정 부분 생략 가능
         정제 과정은 토큰화 이전과 이후에 지속적으로 진행
            하지만 완벽한 정제가 불가능하기 때문에
'amazonhelp']
                    적절한 합의점 선택
                                     nltk 패키지의 Porter Stemmer 사용
```

→ 어간 추출을 통해 arrive -> arriv로 변환

# 2

단어 표현 방법

희소 표현(Sparse representation)

희소 표현(Sparse representation)

벡터 또는 행렬의 값이 대부분 0으로 표현되는 방법 이산적으로 단어의 의미 표현

Book과 Books같은 단어의 변형 형태도 다른 단어로 간주

#### **One-Hot Encoding**

#### One-Hot vector

N개의 단어를 N차원의 벡터로 나타낸 것

John likes to watch movies. Mary likes movies too.

John	1	0	0	0	0	0	0
likes	0	1	0	0	0	0	0
to	0	0	1	0	0	0	0
watch	0	0	0	1	0	0	0
movies	0	0	0	0	1	0	0
mary	0	0	0	0	0	1	0
too	0	0	0	0	0	0	1

단어의 존재 여부에 따라 1 또는 0으로 값이 정해짐

#### **One-Hot Encoding**

#### One-Hot vector

N개의 단어를 N차원의 벡터로 나타낸 것

John likes to watch movies. Mary likes movies too.

John	1	0	0	0	0	0	0	
likes	0	1	0	0	0	0	0	
단어 집	  한에 포호	남되는 단C	거가 많아?	》 <u>직수록</u> 차	원이 기 <u>하</u>	급수적으	로 커짐	
watch	0	- 0	1 1 10 11	_ ' 1' '	_ '0 ' '	0		
movies	0		0	- ' ' ' '메모리 닝	0	0	0	
	0		0	1	0	0	0	

단어의 존재 여부에 따라 1 또는 0으로 값이 정해짐

#### Bag of Words(BoW)

Bag of Words(BoW)

문서 내 단어들에 대해 각 단어의 <mark>등장 빈도를 세는</mark> 방법 단어들의 순서를 전혀 고려하지 않음

John likes to watch movies. Mary likes movies too.



```
{"John":0, "likes":1, "to":2, "watch":3, "movies":4, "Mary":5, "too":6}
```

Word to Frequencies: [1, 2, 1, 1, 2, 1, 1]

단어에 인덱스 부여 후 각 단어의 빈도 수 계산

#### Document-Term Matrix(DTM)

#### Document-Term Matrix(DTM)

문서가 여러 개일 때, 각 문서마다 단어 빈도를 행렬로 나타내는 방법

Document1: John likes to watch movies. Mary likes movies too.

Document2: I think that movie was awful, but John likes that movie.

Document3: John likes Mary.



	John	likes	to	watch	movies	Mary	too	I	think	that	movie	was	awful
Doc1	1	2	1	1	2	1	1	0	0	0	0	0	0
Doc2	1	1	0	0	0	0	0	2	1	2	2	1	1
Doc3	1	1	0	0	0	1	0	0	0	0	0	0	0

위의 예시처럼 문서 별 빈도 수 행렬을 만들 수 있음

#### Document-Term Matrix (DTM)



Document-Term Matrix (DTM)

문서가 여러 개일 때, 리**모TM의 문전점** 행렬로 나타내는 방법

Document1: Johr불용어처럼 자주 나오쁜 단어라면, vies too.

Document2:단어의 중요도에 상관없이 높은 빈모수 가짐

Document3: John likes Mary.



	John	likes	to	watch	movies	Mary	too	I	think	that	movie	was	awful
Doc1	1	2	1 7	<u>5</u> 95	타나	낮은	단어	들어	대혀	H ø	0	0	0
Doc2	1	1 5	드 타-	른《조	大  추	리해		것이	으 오 양	선택	<u>.</u> H 2	1	1
Doc3	1	1	0	0	0	1	0	0	0	0	0	0	0

문서 별 빈도수 행렬

**TF-IDF**(Term Frequency-Inverse Document Frequency)

**TF-IDF**(Term Frequency-Inverse Document Frequency)

단어의 빈도와 역 문서 빈도(문서의 빈도에 특정 식을 취함)를 곱하는 방식



>> DTM 내 각 단어들마다 중요한 정도를 가중치로 줄 수 있음

역 문서 빈도 (IDF)

$$idf(t) = log \frac{n}{1 + df(t)}$$
  $t$  : 문서 내의 단어 수  $n$  : 전체 문서 수

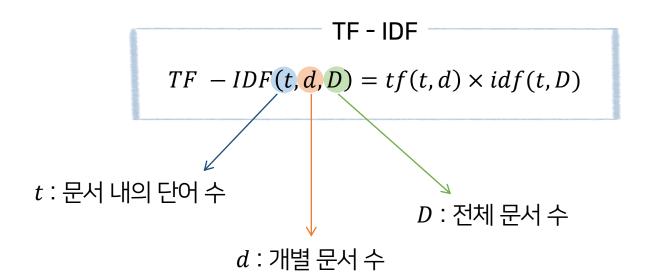
분모가 0이 되지 않기 위해 1을 더함 과도한 가중치를 주는 것을

방지하기 위해 로그를 취함

**TF-IDF**(Term Frequency-Inverse Document Frequency)

**TF-IDF**(Term Frequency-Inverse Document Frequency)

단어의 빈도와 역 문서 빈도(문서의 빈도에 특정 식을 취함)를 곱하는 방식



#### 희소 단어 표현의 문제점

#### 고차원 문제

희소 표현은 전체 단어 집합 크기만큼의 차원을 가지는 벡터 생성 많은 단어가 존재하는 경우, 매우 고차원의 벡터 생성

→ 계산 비용 기하급수적 증가 & 차원의 저주로 효율적 학습 X

#### 공간의 낭비(Sparse Matrix)

희소 표현 방법에서는 <mark>대부분의 값이 0</mark>이 됨 즉, 대부분은 정보를 가지고 있지 않음

→ 텍스트 데이터 정제와 불용어 처리 과정에서 어느 정도 보완 가능 but 여전히 심한 공간의 낭비는 불가피

#### 희소 단어 표현의 문제점

#### 단어 간 의미 관계 부재

One-Hot Encoding & Bow 방법: 단어 존재 유무, 단어 빈도 수를 정수형으로 나타냄임의로 부여한 숫자로, 단어 간의 의미적 관계가 포함되어 있지 않음

모델이 단어 간의 의미적 유사성을 학습하는 데 제약 발생

#### 밀집 표현(Dense Representation)

밀집 표현(Dense Representation)

단어 벡터의 차원을 단어 집합의 크기로 상정하지 않고, 사용자가 설정한 값으로 모든 단어의 벡터 표현의 차원을 맞추는 방식 밀집 표현 과정에서 실수 값을 가지게 됨

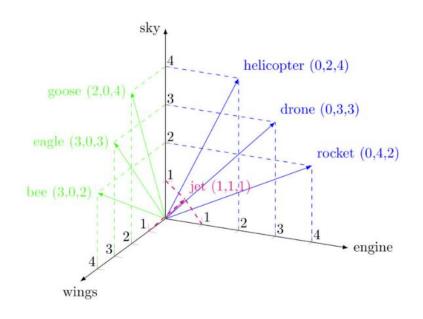
강아지 = [0.2 1.8 1.1 -2.1 1.1 2.8 ... 중략 …]

의미 부여 가능한 128 차원의 벡터로 표현

#### 모델 기반 Word Embedding

#### 모델 기반 Word Embedding

각 단어를 인공 신경망 학습을 통해 벡터화하는 Word Embedding 방법 Word Embedding은 단어를 밀집 벡터의 형태로 표현하는 것

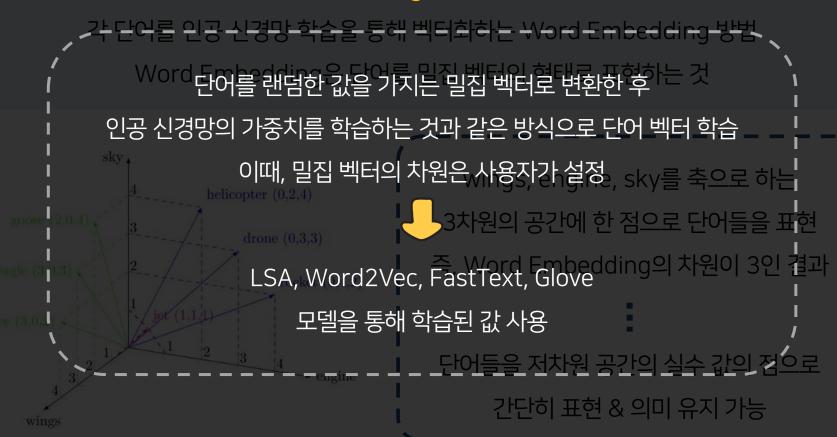


단어들을 저차원 공간의 실수 값의 점으로 간단히 표현 & 의미 유지 가능

#### 모델 기반 Word Embedding



#### 모델 기반 Word Eword Embedding은 어떻게 계산할까?



#### Word2Vec

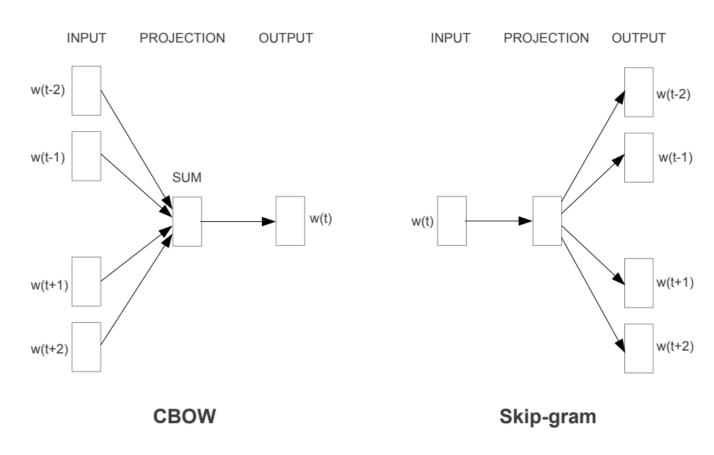
#### Word2Vec

단어를 밀집(dense) 벡터로 표현하는 방법 중 하나로, 임베딩 벡터를 학습하기 위해 주변 단어의 문맥을 활용



자연어 처리에서 단어의 의미를 효과적으로 포착하기 위해 사용 주변 문맥을 어디까지 활용할 것인지 사용자가 지정해야 함

#### Word2Vec



주변 문맥 학습 방법에 따라 CBOW 모델과 Skip-gram 모델로 나뉨

#### Word2Vec

CBOW: 주변 단어 문맥으로 중심 단어를 예측하는 방식

$$\frac{1}{T} \sum_{t=1}^{T} \sum_{-c \le j \le c, j \ne 0} \log p(w_t | w_{t+j})$$

각각의 확률을 최대로 하는 단어를 찾도록 학습

Skip-gram: 중심 단어에서 주변 단어를 예측하는 방식

$$\frac{1}{T} \sum_{t=1}^{T} \sum_{-c \le j \le c, j \ne 0} \log p(w_{t+j}|w_t)$$

일반적으로 CBOW 모델에 비해, Skip-gram 모델로 학습한 임베딩의 성능이 더 좋은 것으로 알려짐

#### Word2Vec



CBOW : 주변 단어 문맥으로 중심 단어를 예측하는 방식 , 학습을 위한 데이터셋은 무엇을 사용하는 것이 좋을까?  $\overline{\log p(w_t|w_{t+j})}$ 각각의 확률을 최대로 하는 BoW 방법을 응용해 단어를 찾도록 학습 중심 단어에서 주변 단어를 예측하는 방식 같이 나오는 단어들의 빈도 수를 저장한 행렬인 co-occurrence matrix를 데이터셋으로 하여 학습 -c≤j≤c,j≠0 일반적으로 CBOW 모델에 비해

Skip-gram 모델로 학습한 임베딩의 성능이 더 좋은 것으로 알려짐

#### Co-occurrence matrix 예시

Document1: Apples are green and red.

Document2: Red apples are sweet.

Document3: Green oranges are sour.

Document4: Monkeys eat red apples.

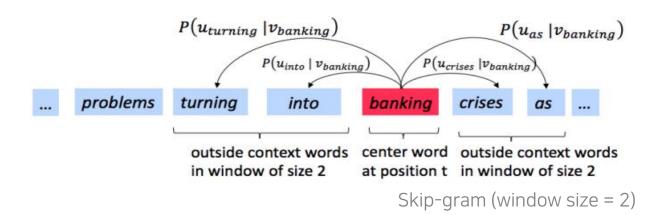


-	apples	are	green	and	red	sweet	oranges	sour	monkeys	eat
apples	3	2	1	1	3	1	0	0	1	1
are	2	3	2	1	2	1	1	1	0	0
green	1	2	2	1	1	0	1	1	0	0
and	1	1	1	1	1	0	0	0	0	0
red	3	2	1	1	3	1	0	0	1	1
sweet	1	1	0	0	1	1	0	0	0	0
oranges	0	1	1	0	0	0	1	1	0	0
sour	0	1	1	0	0	0	1	1	0	0
monkeys	1	0	0	0	1	0	0	0	1	1
eat	1	0	0	0	1	0	0	0	1	1

### Skip-gram

Skip-gram

중심 단어를 기준으로 주변의 단어 예측 앞 뒤로 볼 단어의 범위를 윈도우라고 함



### Skip-gram

#### 주변 단어들의 확률 분포

$$P(O = o | C = c) = \frac{\exp(\boldsymbol{u_0^T v_c})}{\sum_{w \in Vocab} \exp(\boldsymbol{u_W^T v_c})}$$

$$\mathbf{J}(\mathbf{v_c}, \mathbf{o}, \mathbf{U}) = -logP(O = o | C = c)$$

이 확<del>률을</del> 최대화하는 방향으로 학습

-1을 곱하여 손실을 양의 값으로 바꾸고, 최대화 하는 방향으로 학습 가능

SGNS(Skip-gram with Negative Sampling)

SGNS(Skip-gram with Negative Sampling)

좀 더 빠른 계산을 위해 Skip-gram에 Negative Sampling 기법을 추가한 모델

### **Negative Sampling**

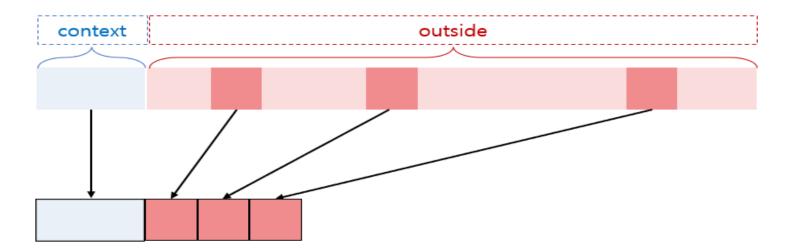
Word2Vec이 학습 과정에서 전체 단어 집합이 아니라 일부 단어 집합에만 집중할 수 있도록 하는 방법

### **Negative Sampling**

$$P(O = o | C = c) = \frac{\exp(\mathbf{u_0^T v_c})}{\sum_{w \in Vocab} \exp(\mathbf{u_w^T v_c})}$$
 중심 벡터의 전체 벡터에 대한 dot product 계산

일반적인 경우, 전체 단어 집합에 대해 계산하므로 단어 집합이 커질수록 연산 양이 훨씬 더 많아짐

#### **Negative Sampling**



전체 단어 집합에는 속하지만 주변 단어에는 속하지 않는 단어 중,

**랜덤하게 샘플링**해 부분 단어 집합을 만듦



전체 집합 대신 부분 집합에 대해 계산

SGNS(Skip-gram with Negative Sampling)

소실 함수 
$$\mathbf{J}(\mathbf{v_c}, \mathbf{o}, \mathbf{U}) = -log(\sigma(\mathbf{u_0^T v_c}) - \sum_{s=1}^K \log(\sigma(\mathbf{-u_{W_s}^T v_c}))$$

내적 계산 시 -1을 곱하여 outside vector를 포함하지 않는 방향으로 학습함

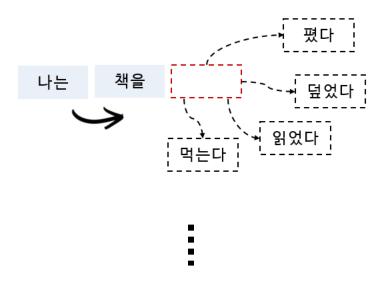
# 3

언어 모델링

### 언어 모델링

언어 모델링(Language Modeling)

언어라는 현상을 모델링하고자 문장(단어 시퀀스)에 확률을 할당하는 모델



단어 시퀀스에 확률을 할당하기 위해 이전 단어들로부터 다음 단어 예측

언어 모델링(Language Modeling)

언어라는 현상을 모델링하고자 <mark>문장(단어 시퀀스)에 확률을 할당</mark>하는 모델



즉, 이전 단어들을 입력으로 <mark>조건부 확률</mark>을 생성하는 것

$$P(Y) = P(y_1, ..., y_t) = \prod_{t=1}^{\tau} P(y_t | y_{t-1}, ..., y_1)$$

통계적 언어 모델(SML)

인공 신경망 기반 모델(NML)

# 통계적 언어 모델(SML)

#### 카운트 기반 접근 모델

 $x^{(t+1)}$ 은 오직 이전 t개의 항에 의존한다는 <mark>마르코프 가정</mark>에 기반함

#### N-gram

전체를 보는 것이 아닌 이전의 (N-1)개의 단어만 보는 방법

### 통계적 언어 모델(SML) – 카운트 기반 접근 모델

카운트 기반 접근 모델

 $x^{(t+1)}$ 은 오직 이전 t개의 항에 의존한다는 <mark>마르코프 가정</mark>에 기반함

N-gram

전체를 보는 것이 아닌 이전의 (N-1)개의 단어만 보는 방법

Ex) P (나는 밥을 먹고 다시 학교로 가서 공부를 했다) = P(나는) X P(밥을|나는) X …

X P ( 했다 | 나는 밥을 먹고 다시 학교로 가서 공부를 )



count ( 나는 밥을 먹고 다시 학교로 가서 공부를 했다 ) count ( 나는 밥을 먹고 다시 학교로 가서 공부를 )

통계적 언어 모델(SML) – 카운트 기반 걸근 모델

카운트 기반 접근 모델

χ<sup>(t+1)</sup>은 오직**만약 카운트가 이이 된다면?"** 

 $\stackrel{\text{\tiny 44}}{\longleftarrow}$  분모가 0인 경우: 분모에 매우 작은 수( $\epsilon$ )를 더함

Ex) 👑 분자가 0인 경우: 더할 수 없으므로,했다)

훈련 데이터에서 실제 count가 없다면 확률은 0

X<mark>I</mark>P ( 했다 | 나는 밥을 먹고 <u>다</u>시 학교로 가서 공부를 )

데이터가 매우 드물게 분포되어 있어 모델 학습에 어려움을 겪는 희소 문제 발생

count ( 나는 밥을 먹고 다시 학교로 가서 공부를 )

### 통계적 언어 모델(SML) - N-gram

카운트 기반 접근 모델

 $x^{(t+1)}$ 은 오직 이전 t개의 항에 의존한다는 **마르코프 가정**에 기반함

N-gram

전체를 보는 것이 아닌 이전의 (N-1)개의 단어만 보는 방법

나는 밥을 먹고 다시 학교로 가서 공부를 \_\_\_\_

무시!

4개의 단어만 고려!

N-gram(n = 5)의 예시

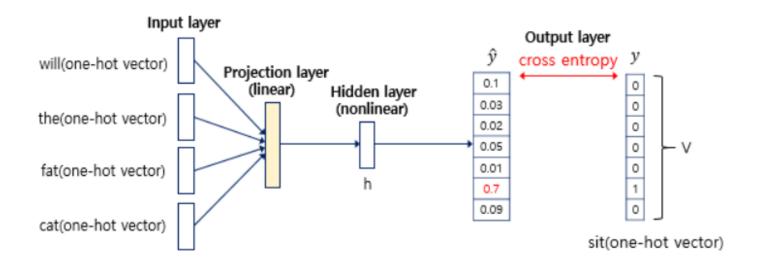


이 역시 희소 문제를 피할 수 없었음

# 인공 신경망 기반 언어 모델(NML)

일반적인 인공 신경망(Feed-forward)

밀집 표현인 임베딩을 학습함으로써 해결함



인공 신경망 기반 언어 모델(NML)



카운트 기반 접근 모델

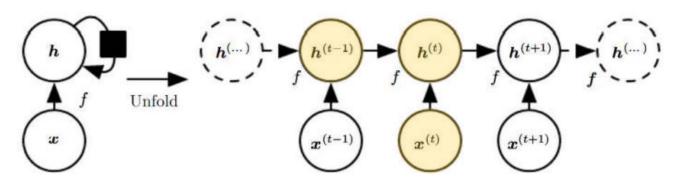
x(t+1)은Feed-forward 신경망의 문제점 기반함

N-gram 고정된 길이의 단어 시퀀스만을 입력으로 받기 때문에 자연어 처리를 위한 방법으로는 한계점이 존재함 일반적인 인공 신경망 (Feed-forwa 가변적인 길이를 처리할 수 있는 새로운 모델의 필요성 학습함으로써 해결함 RNN 등장

#### Vanilla RNN



모델의 Output이 Input뿐만 아니라 이전의 Output에도 영향을 받음



왼: 모델의 순환성 / 오: 시점에 따른 학습으로 펼친 모습

t 시점에는 t번째 인풋, t-1시점의 아웃풋을 입력으로 받음

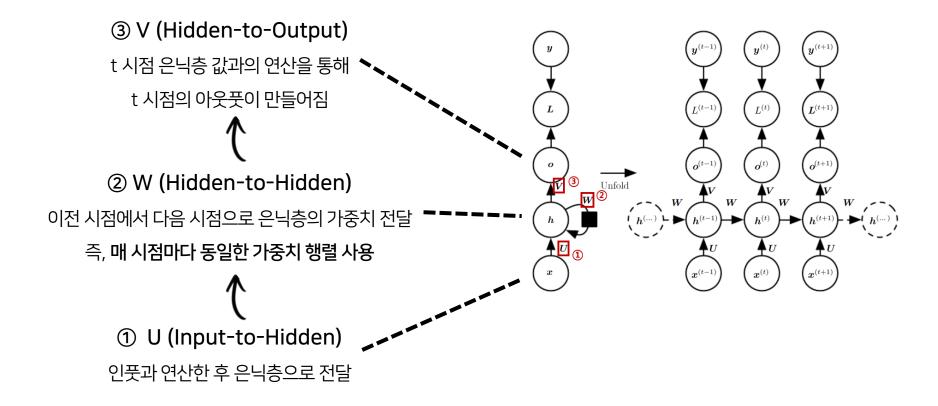


기존 인공 신경망과 달리 가변적 길이의 단어 시퀀스를 입력으로 받음

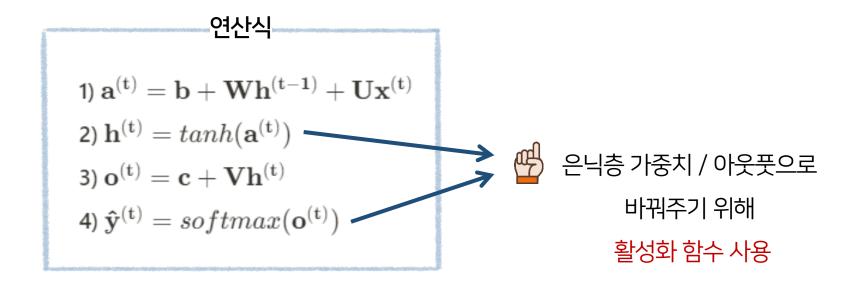
### RNN의 연산



#### 3개의 가중치 행렬



#### RNN의 연산



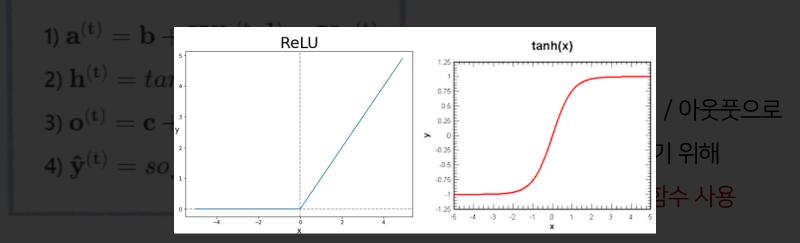
∰ 매개변수들이 시점 t에 의존하고 있지 않음

즉, 순전파 과정에서 매개변수는 모든 시점에서 공유되며 변하지 않음

RNN의 연산

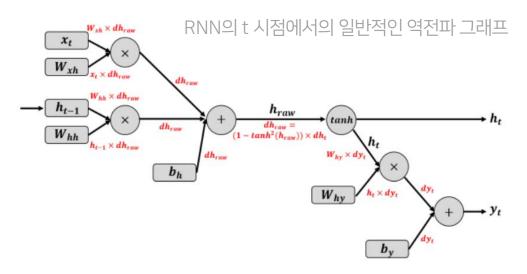


# RNN에서 ReLU 함수를 사용하지 않는 이유



ReLU는 절댓값 1 이상의 값을 가질 수 있음 <u>순환적인 RNN의 특성 상 발산의 가능성을 만들어주기에 사용하지 않음</u>

#### RNN 역전파



RNN은 시점에 따라 다른 출력이 존재

순전파 과정이 끝나 모든 출력이 만들어진 다음 마지막 출력부터 역전파가 일어나야 함



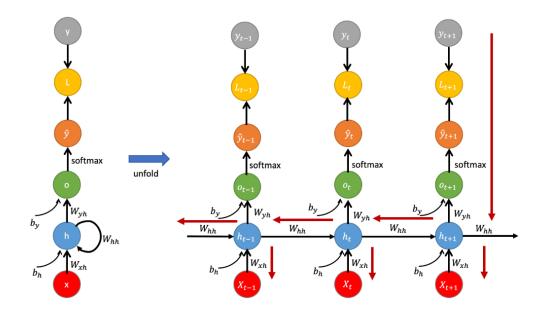
BPTT 사용

₩ 각 시점마다 은닉층 상태를 메모리에 저장하고 있어야 함

**BPTT**(Back Propagation through Time)

**BPTT** (Back Propagation through Time)

기울기가 매 시점마다 계산되어 손실 값 또한 매 시점 계산되는 형태의 역전파



모델을 시점 단위로 펼쳐서 마지막 시점부터 첫 번째 시점까지 업데이트

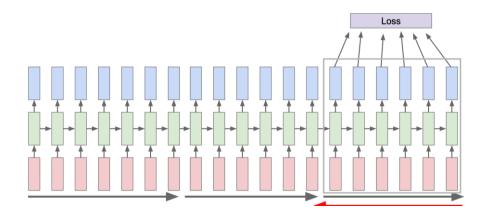
#### Truncated BPTT

Sequential Data의 길이가 길어지면, BPTT는 더 오래 걸림 값이 1보다 크게/작게 되면 첫 시점으로 갈수록 미분 값이 커짐/작아짐 즉, 과거 시점에 대한 예측이 거의 불가능함



#### **Truncated BPTT**

특정 시점 단위로 끊어서 기울기를 계산하는 형태의 역전파



### RNN의 한계점 - 1) 매우 긴 학습 시간

출력 결과가 이전의 계산 결과에 의존함



즉, 이전의 계산 결과가 끝나야지만 다음 단어 시퀀스에 계산이 가능하므로 loop를 통해 계산할 수밖에 없음

Ex) 나는 밥을 먹고 다시 학교로 가서 공부를 했다

단어 8개 = 8번의 loop



매우 느린 학습 시간 초래

RNN의 한계점



# BPTT를 통해 미분값을 계산할 때의 문제점

시퀀스가 길어짐에 따라

기울기 < 1

기울기 > 1

: 매우 작은 기울기 값이 됨

: 훨씬 더 큰 기울기 값이 됨

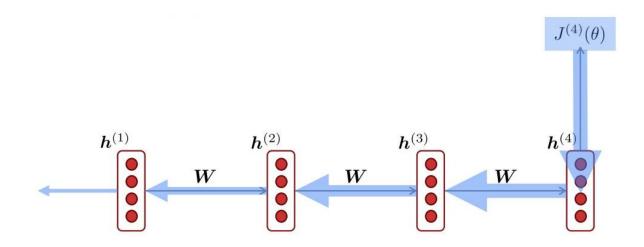
시점이 멀어짐에 따라 정보가 제대로 반영되지 못 함

즉, 장기 의존성 문제가 발생함

#### RNN의 한계점 - 2) 장기 의존성 문제

장기 의존성 문제

은닉층의 과거의 정보가 마지막까지 제대로 전달되지 못 하는 현상



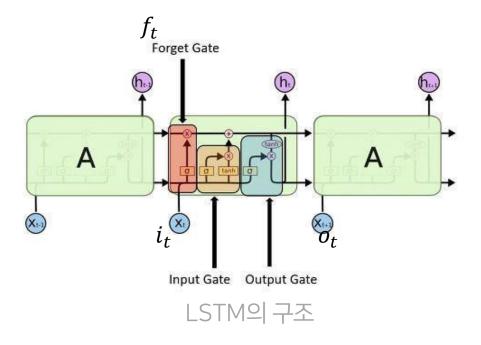
가까운 시점의 데이터 위주로 학습이 진행됨

# LSTM (Long Short-Term Memory)

**LSTM** 

단기 기억을  $h_t$ (hidden state), 장기 기억을  $c_t$ (cell state)

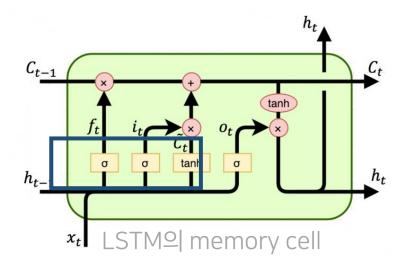
각각 담당해 두 개의 hidden vector를 갖는 모델



### LSTM (Long Short-Term Memory)

**LSTM** 

단기 기억을  $h_t$ (hidden state), 장기 기억을  $c_t$ (cell state) 각각 담당해 두 개의 hidden vector를 갖는 모델



3개의 게이트를 통해 어떤 정보를 읽고, 지우고, 쓸 것인지 결정함

### LSTM의 연산 – 1) 게이트 연산

$$f^{(t)} = \sigma(W_f h^{(t-1)} + U_f x^{(t)} + b_f)$$

$$i^{(t)} = \sigma(W_i h^{(t-1)} + U_i x^{(t)} + b_i)$$

$$o^{(t)} = \sigma(W_o h^{(t-1)} + U_o x^{(t)} + b_o)$$





**RNN의 연산과 동일하며** 

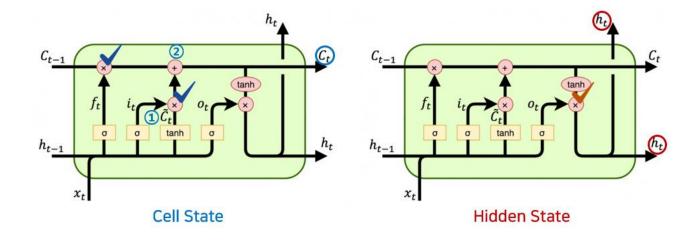
가중치는 시점과 무관하게 존재함



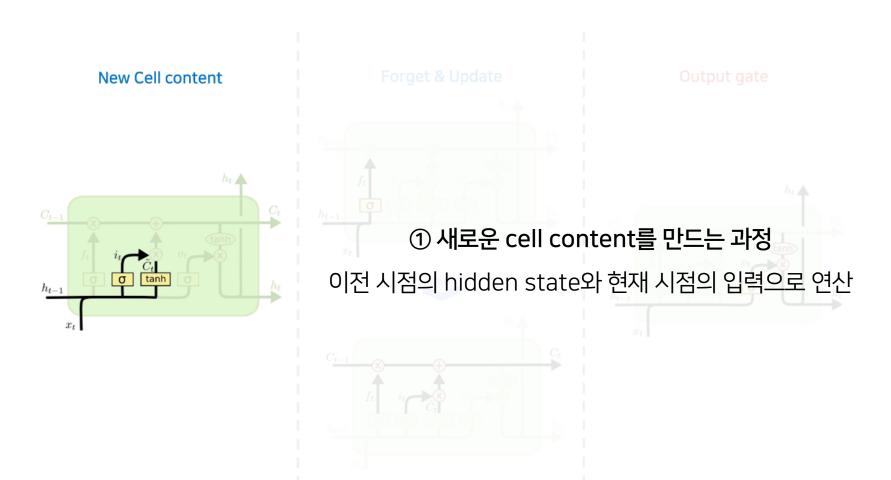
🖺 0과 1 사이의 값을 갖는 Sigmoid 함수를 적용함으로써

정보를 얼마나 반영하는지에 대한 가중치로 작용함.

# LSTM의 연산 - 2) 셀 상태 & 은닉 상태 연산



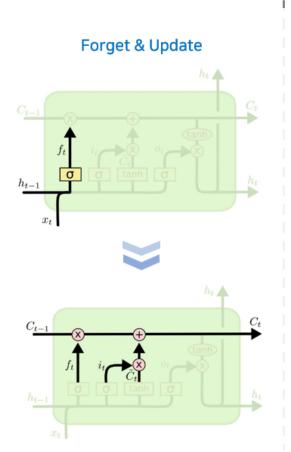
### LSTM의 연산 - 2) 셀 상태 & 은닉 상태 연산



### LSTM의 연산 - 2) 셀 상태 & 은닉 상태 연산

② 현재 시점의 cell state를 만드는 과정





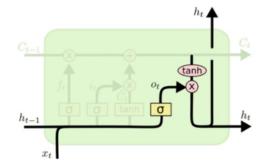
As a

이전 상태의 정보를 얼마나 잊고 현재 상태의 정보를 얼마나 반영할지 결정

### LSTM의 연산 - 2) 셀 상태 & 은닉 상태 연산

③ 새로운 hidden state를 만들어내는 과정 Output gate의 결과물과 현재 시점의 cell state를 곱함

**Output gate** 



LSTM의 연산





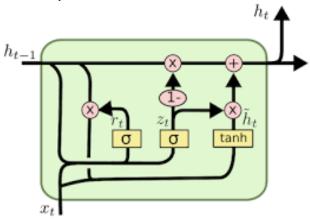
### **GRU (Gate Recurrent Unit)**



LSTM: 3개의 게이트(input, forget, output)



GRU: 2개의 게이트(update, reset)로 줄어들어 구조가 간단해짐



일반적으로 GRU의 학습 속도가 더 빠르다고 알려져 있으나, 성능은 비슷함

#### Other RNNs

#### **Bidirectional RNN**

이전 시점의 정보만을 활용하는 Vanilla RNN과 달리 양방향으로 RNN을 학습시키는 구조

#### Multi-layer RNN

Neural Network를 여러 층으로 쌓아 RNN에 더 복잡한 구조를 학습 시키는 방법

#### Other RNNs

#### **Bidirectional RNN**

이전 시점의 정보만을 활용하는 Vanilla RNN과 달리 양방향으로 RNN을 학습시키는 구조

RNN 대신 LSTM, GRU로 대체 가능

The movie was terribly exciting!

exciting이 terribly의 의미를 <mark>부정에서 긍정</mark>으로 전환



단방향 RNN의 경우 문맥에 따른 의미 변화를 파악하기 어려워 terribly를 부정적 의미로만 파악

#### Other RNNs

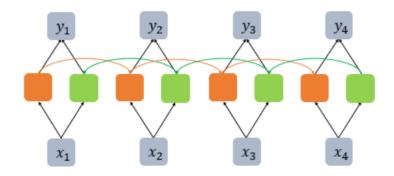
#### **Bidirectional RNN**

이전 시점의 정보만을 활용하는 Vanilla RNN과 달리 양방향으로 RNN을 학습시키는 구조

RNN 대신 LSTM, GRU로 대체 가능

The movie was terribly exciting!

exciting이 terribly의 의미를 부정에서 긍정으로 전환



#### Other RNNs

#### **Bidirectional RNN**

이전 시점의 정보만을 활용하는 Vanilla RNN과 달리 양방향으로 RNN을 학습시키는 구조

RNN 대신 LSTM, GRU로 대체 가능

#### Forward RNN

$$\vec{h}^{(t)} = RNN_{FW}(\vec{h}^{(t-1)}, x^{(t)})$$

#### **Backward RNN**

$$\widetilde{h}^{(t)} = RNN_{BW}(\widetilde{h}^{(t-1)}, x^{(t)})$$

#### Concatenated RNN

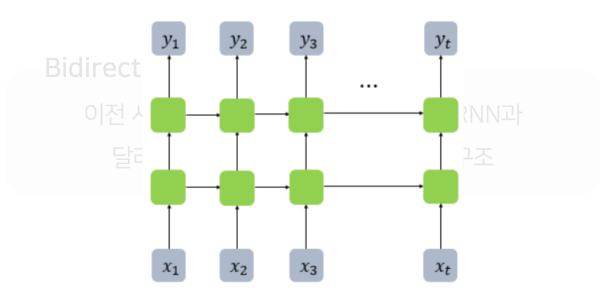
$$h^{(t)} = [\overrightarrow{h}^{(t)}; \overleftarrow{h}^{(t)}]$$

### Other RNNs

### Multi-layer RNN

Neural Network를 여러 층으로 쌓아 RNN에 더 복잡한 구조를 학습 시키는 방법

RNN 대신 LSTM, GRU로 대체 가능



Other RNNs



# 언어 모델링의 성능 평가는 어떻게 할까?

Neural Network를 여러 층으로 쌓아

언어 모델링에 다양한 모델이 존재

해당 모델들을 비교하기 위해서 평가지표가 필요

회귀문제에서 사용하는 RMSE나

분류문제에서 사용하는 정확도와 같은

평가 방법인 PPL 사용

### Perplexity(PPL)

#### Perplexity(PPL)

언어 모델링에서 모델 자체적으로 자신의 성능을 수치화하기 위해 사용하는 지표 문장의 길이로 정규화 된 문장 확률의 역수



'헷갈리는 정도'를 의미 PPL이 낮을수록 언어 모델의 성능이 좋음

### Perplexity(PPL)

PPL

$$PPL(W) = P(w_1, w_2, ..., w_N)^{-\frac{1}{N}} = \sqrt[N]{\frac{1}{P(w_1, w_2, ..., w_N)}} = \sqrt[N]{\frac{1}{\pi_{i=1}^N P(w_i | w_1, w_2, ..., w_N)}}$$



N-gram 적용 (N = 5)

$$PPL(W) = \sqrt[N]{\frac{1}{\pi_{i=5}^{N} P(w_i|w_{i-4}, w_{i-3}, w_{i-2}, w_{i-1})}}$$

### Perplexity(PPL)

PPL

문장 W의 길이가 N일 때,

$$PPL(W) = P(w_1, w_2, ..., w_N)^{-\frac{1}{N}} = \sqrt[N]{\frac{1}{P(w_1, w_2, ..., w_N)}} = \sqrt[N]{\frac{1}{\pi_{i=1}^N P(w_i | w_1, w_2, ..., w_N)}}$$



#### 문장 내에서 평균적으로 몇 개의 선택지를 가지고 고민하는지를 의미

$$PPL(W) = \sqrt[N]{\frac{1}{\pi_{i=5}^{N} P(w_i|w_{i-4}, w_{i-3}, w_{i-2}, w_{i-1})}}$$

### PPL 계산 예시

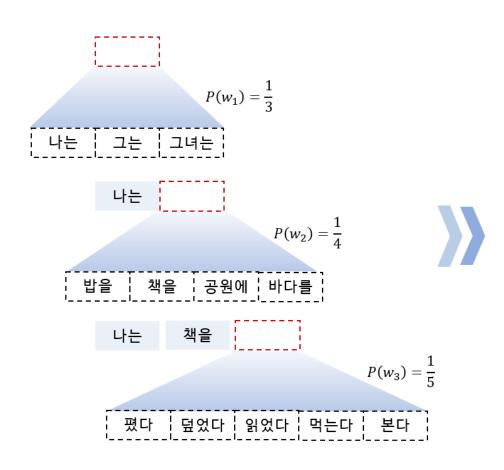
$$A : PPL = 20 / B : PPL = 30$$

문장의 모든 시점마다 모델 A는 20개, 모델 B는 30개의 선택지 고민



모델 A가 B에 비해 덜 헷갈려 하므로 A가 더 성능이 좋음

### PPL 계산 예시



실제 정답: "나는 책을 읽었다"

PPL = 4

# 4

## 기계 번역

### 기계 번역

#### 기계 번역

어떤 문장 X를 다른 문장 Y로 번역하는 것

언어 모델링과 더불어 자연어 처리 분야에서 대표적인 과제 중 하나!

X: 나는 NLP 팀원들이 있어서 안심이 돼

Y: I am relieved to have NLP team members

기계 번역의 주요 이슈



대응되는 단어의 부재

대응되는 단어가 없는 경우 **→** 한국어에는 조사가 있지만, 영어에는 그렇지 않음

나는 NLP 팀원들이 있어서 안심이 돼!

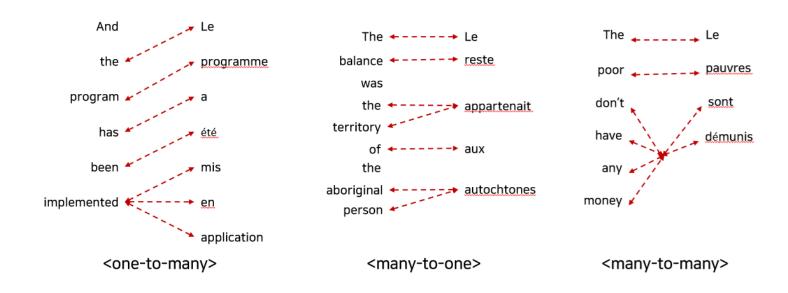
는, 이 등과 같은 조사는 영어에 존재하지 않음

#### 기계 번역의 주요 이슈



정렬(Alignment)의 문제점

어떤 단어는 하나의 단어에만 대응 (one-to-one) 대응 관계에서 정렬에 대한 문제점 → 어떤 단어는 여러 개에 대응 (one-to-many) 여러 개의 단어가 여러 개에 대응 (many-to-many)



기계 번역의 종류

기계 번역에는 통계적 기계 번역과 인공 신경망 기계 번역이 있음

통계적 기계 번역 (SMT) 인공 신경망 기계 번역 (NMT)

**SMT**(Statistical Machine Translation)

**SMT**(Statistical Machine Translation)

데이터로부터 확률적 모델을 학습하는

통계적 접근 방법

"프랑스어를 한국어로 번역해보자!"

프랑스어 문장 x

가장 좋은 한국어 번역 문장
y

**SMT**(Statistical Machine Translation)

이때, 베이지안 법칙에 기반하여 두 가지 문제로 쪼갤 수 있음  $argmax_y P(y|x) = argmax_y (P(x|y)P(y))$ 

P(x|y) : 한국어 문장이 주어졌을 때, 프랑스어 문장으로 번역하는 방법에 대해 학습하는 번역 모델(MT)

 P(y)
 : 좋은(유창한) 한국어 문장을 만들어내는 것을 목표로 하는

 언어 모델(LM)

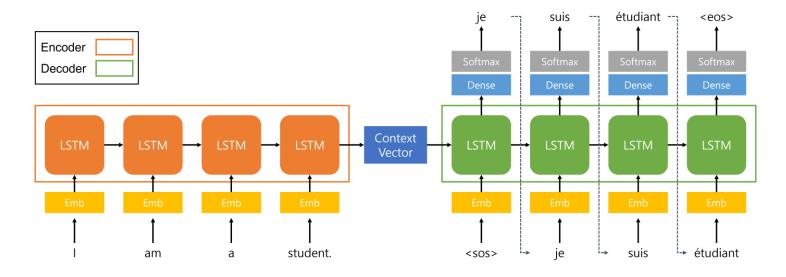
그러나, SMT는 큰 성과를 거두지 못 함.

→ 2014년, 하나의 end-to-end 모델로 기계번역을 하는 NMT(Neural Machine Translation)이 등장

NMT(Neural Machine Translation)

NMT(Neural Machine Translation)

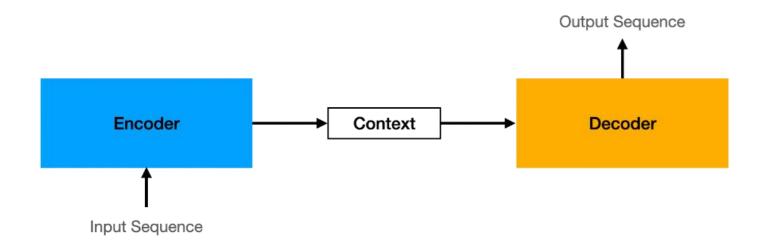
인공 신경망을 사용하는 기계 번역 접근 방법



하나의 end-to-end 모델로 기계 번역이 가능함

: 입력에서 출력까지 신경망으로 한 번에 처리

#### **Encoder-Decoder**



원래의 문장을 압축하여

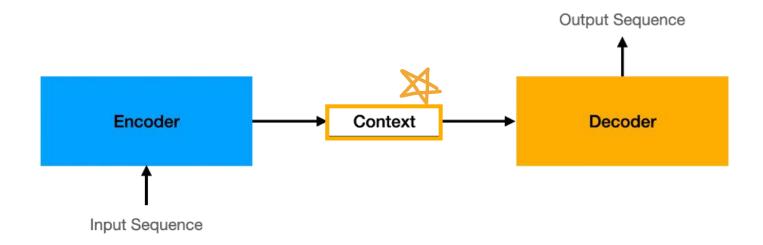
<mark>하나의 벡터</mark>로 만듦

만들어진 벡터를 바탕으로 새로운 문장으로 복원



context vector

#### **Encoder-Decoder**



Context vector는 **Encoder**의 마지막 hidden state 값이 됨



이 context vector를 받아서

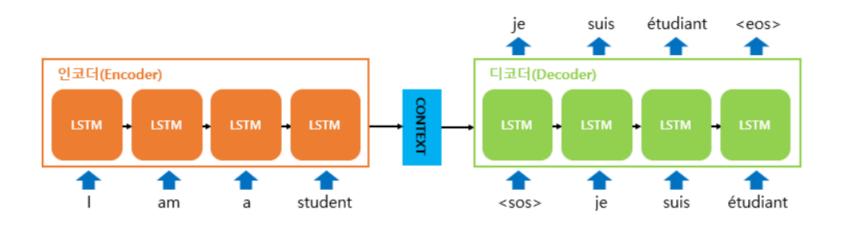
Decoder의 첫 번째 hidden state로 사용

### seq2seq

#### sequence-to-sequence model

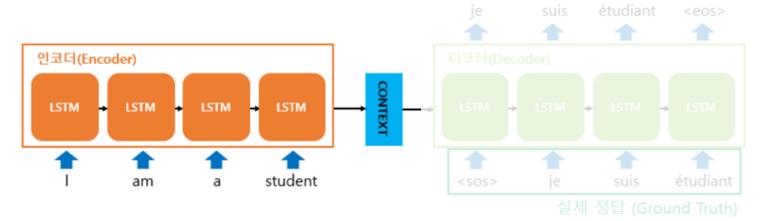
NMT에서 대표적으로 사용되는 Encoder-Decoder 모델 Encoder와 Decoder에 각각 RNN 계열의 모델을 사용

특히 내부 cell로 LSTM 을 사용하는 것이 학습 성능이 좋음



#### seq2seq에서 Encoder-Decoder의 학습

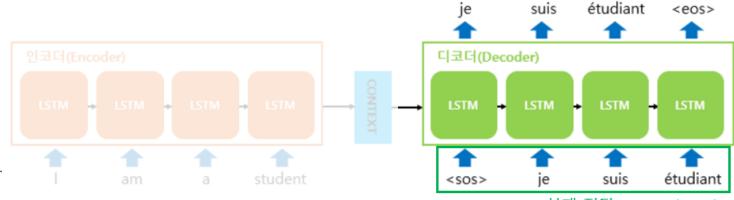
#### Encoder 학습



입력 문장을 받아서 하나의 context vector로 압축하고 입력 문장의 문맥을 학습함

#### seq2seq에서 Encoder-Decoder의 학습

#### Decoder 학습



실제 정답 (Ground Truth)

Decoder는 context vector와 새로운 입력에 의존함

테스트 과정 - 학습 과정에서의 작동 방식이 다름

#### seq2seq에서 Encoder-Decoder의 학습

#### Decoder 학습



<sos>\*: Start-Of-String, 문장의 시작을 알리는 토큰

첫 번째 셀은 context vector와 <sos>\*를 바탕으로 je를 생성 두 번째 셀은 첫 번째 셀의 예측인 je를 바탕으로 suis를 생성

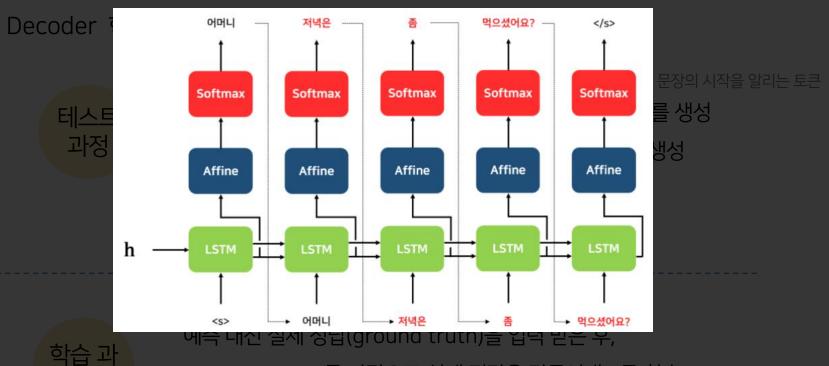
:

이 과정을 문장이 끝날 때까지 반복!



예측 대신 실제 정답(ground truth)을 입력 받은 후, context vector를 바탕으로 실제 정답을 만들어내도록 학습하는 교사강요 적용

### seq2seq에 Pecceder 학습에서 교사 강요를 사용하는 이유



context vector를 바탕으로 실제 정답을 만들어내도록 학습 Teacher Forcing을 **사용하지 않을 경우** 

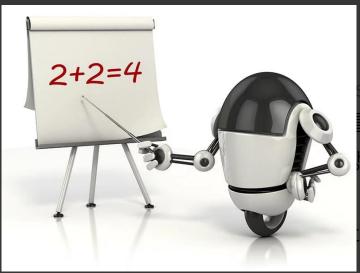
: t번째 예측을 만들어내기 위해 t-1번째 디코더의 예측값을 사용

만약 예측이 올바르지 않다면, <mark>이전 디코더의 예측 값</mark>을 사용하는 것은 큰 문제

## seq2seq에 Pecceder 학습에서 교사 강요를 사용하는 이유

Decoder 학습

테스트 과정



를 바탕으로 je를 생성

으로 suis를 생성

ll까지 반복 !

장점

예측 대신 실제 정답(ground truth)을 입력 받은 후.



Teacher Forcing은 학습 과정에서 실제 정답을 입력 받음으로써 **더 빨리 학습 가능** 

context vector를 바탕으로 실제 정답을 만들어내도록 학습하는 교 학습과 추론 사이의 차이가 생기는

노출 편향 문제(Exposure Bias Problem)가

존재하나 성능에 큰 영향을 미치지 않음

### seq2seq의 한계점



고정된 길이의 context vector

입력 문장의 모든 정보를 하나의 벡터에 저장 → 정보의 손실 발생 길이가 매우 길다면 고정된 길이의 context vector에 이를 다 반영하는 것은 어려움



#### 기울기 소실 문제

이전의 은닉 상태 $(h_{t-1})$ 와 현재의 입력 $(x_t)$ 을 바탕으로 현재의 출력을 예측

→ 시간의 흐름에 따른 연속적인 구조로 기울기 소실 문제 존재

#### **BLEU Score**

#### Bilingual Evaluation Understudy

기계 번역 평가 방법으로 "실제 전문가의 번역과 가까울수록 더 좋은 번역이다"라는 내용을 함의

기계 번역 결과 Candidate Translation



사람이 직접 번역한 결과 Reference Translation

#### **BLEU Score**

#### BLEU Score의 장점

- ① 빠르게 점수를 구할 수 있고 계산 비용이 저렴함
- ② 결과를 비교적 이해하기 쉬움
- ③ 언어에 구애 받지 않고 사용할 수 있음
- ④ 사람의 번역과 크게 연관시킬 수 있음
- ⑤ 폭넓게 적용 가능

**BLEU Score** 

$$BLEU = BP \times \exp(\sum_{n=1}^{N} w_n \log p_n)$$

: Brevity Penalty, BP

만들어낸 문장의 길이가 짧을 경우 페널티를 줌

 $\exp(\sum_{n=1}^N w_n \log p_n)$  : 보정된 N-gram 정밀도 기하평균

#### **BLEU Score**

$$BLEU = BP \times \exp(\sum_{n=1}^{N} w_n \log p_n)$$

BP : Brevity Penalty, 만들어낸 문장의 길이가 짧을 경우 페널티를 줌

$$\exp(\sum_{n=1}^N w_n \log p_n)$$
 : 보정된 N-gram 정밀도 기하평균

해당 식이 어떻게 나오게 되는지 알아보자!

### Unigram 정밀도

Candidate 1: It is a guide to action which ensures that the military always oneys the common (Machine Translation)

Reference 2: It is the guiding principle 세계의 Reference Translation (Human Translation) d of the Party

Candidate 
$$1 - (P = \frac{17}{18})$$

C1	It	is	a	guide	to	action	which	ensures	that	the	military	always	obeys	the	commands	of	the	party
R1	0	0	0	0	0	0	Х	0	0	0	0	Х	Х	0	0	Х	0	0
R2	0	0	χ	Х	Х	Х	0	Х	Х	0	0	0	Х	0	0	0	0	0
R3	0	0	Χ	0	0	Х	Х	Х	Х	0	Х	0	Х	0	Х	0	0	0

Candidate 
$$2 - (P = \frac{8}{14})$$

C2	It	is	to	insure	the	troops	forever	hearing	the	activity	guidebook	that	party	direct
R1	0	0	0	Х	0	Х	0	Х	0	Х	Х	0	0	Х
R2	0	0	Х	Х	0	Х	Х	Х	0	Х	Х	Х	0	Х
R3	0	0	0	Х	0	Х	Х	Х	0	Х	Х	Х	0	Х

이전 단어를 고려하지 않고 현재 단어 하나만을 보기 때문에 단어를 단순히 카운트하는 방법

### Unigram 정밀도의 문제점



중복된 단어 처리

Machine Translation: the the the the the the

Human Translation: ① The cat is on the mat

② There is a cat on the mat

기계 번역이 the로만 이루어져 있음

 $\rightarrow$  단어를 단순히 카운트하는 Unigram 방식에서는 결과가  $\frac{7}{7}$  = 1이 나옴

즉, 좋은 번역이 아니어도 좋은 결과가 나올 수 있음

### Unigram 정밀도의 문제점



순서를 고려하지 않음

Unigram 방식은 단어 각각에 대해 독립적으로 카운트하는 방법 그러나, 자연어는 기본적으로 <mark>순차 데이터</mark>

→ 순서를 고려해줄 방법 필요!



보정된 N-gram 정밀도 (Modified N-gram precision)

### 보정된 N-gram 정밀도

중복된 단어 처리 문제 해결

Machine Translation: the the the the the the

Human Translation: 1 The cat is on the mat

② There is a cat on the mat

참조 문장에서는 the가 2번 밖에 안 나왔는데 이를 반영하지 못함

 $Count_{clip} = min(Count, \frac{Max\_Ref\_Count}{})$ 

참조 문장에서 동일한 단어 등장 횟수(최대 참조 횟수)

단순 카운트와 최대 참조 횟수를 동시에 고려함

→ 정밀도는  $\frac{2}{7}$ 가 나옴

### 보정된 N-gram 정밀도

순서 고려

Candidate 
$$1 - (P = \frac{10}{17})$$

C1	It is	is a	a guide	guide to	to action	action which	which ensures	ensures that	that the
R1	0	0	0	0	0	Х	Х	0	0
R2	0	Х	Х	Х	Х	Х	Х	Х	Х
R3	0	Х	Х	Х	Х	Х	Х	Х	Х
C1	the military	military always	always obeys	obeys the	the commands	commands of	of the	the party	
R1	0	Х	Х	Х	Х	Х	Х	Х	
R2	0	Х	Х	Х	Х	Х	0	0	
R3	Х	Х	Х	Х	Х	Х	0	0	

Candidate 
$$2 - (P = \frac{1}{13})$$

C2	it is	is to	to insure	insure the	the troops	troops forever	forever hearing
R1	0	Х	X	Х	Х	Х	Х
R2	0	Х	Х	Х	Х	Х	Х
R3	0	Х	X	Х	Х	Х	Х
C2	hearing the	the activity	activity guidebook	guidebook that	that party	party direct	
R1	Х	Х	Х	Х	Х	Х	
R2	х х		Х	Х	Х	Х	
R3	3 X X		Х	Х	Х	Х	

이전 단어를 포함해 순서를 고려하는 N-gram으로 Unigram의 단점 보완

### Brevity Penalty (BP)

**Brevity Penalty** 

짧은 문장에 대해 페널티를 줌

아래와 같은 Reference 문장과 Candidate 문장에 대해 보정된 2-gram 정밀도로 계산해보자.

Machine Translation: 1 That cat sat on the soft mat.

② Black and white cat sat on mat in corner of room.

Human Translation: The black and white that cat sat peacefully on the soft mat in the corner of the room

#### **Brevity Penalty (BP)**

Candidate 
$$1 - (P = \frac{5}{6})$$

C1	that cat	cat sat	sat on	on the	the soft	soft mat	
R	0	0	Х	0	0	0	

That cat sat on the soft mat.

Candidate 
$$2 - (P = \frac{5}{10})$$

C2	black and	and white	white cat	cat sat	sat on	on mat	mat in	in corner	corner of	of room
R	0	0	Х	0	Х	Х	0	Х	0	Х

Black and white cat sat on mat in corner of room.

Candidate 1은 많은 부분이 생략되어 좋은 문장은 아니지만, 2-gram 정밀도로 계산하면 Candidate 2보다 좋다고 판단

→ 이를 완화하기 위해 BP factor 부여

### **Brevity Penalty (BP)**

$$BP = \begin{cases} 1 & \text{if } c > r \\ e^{(1-r/c)} & \text{if } c \le r \end{cases}$$

Machine Translation: (1) That cat sat on the soft mat.

② Black and white cat sat on mat in corner of room.

Human Translation: The black and white that cat sat peacefully on the

soft mat in the corner of the room.

$$\rightarrow BP = 0.2077$$

$$→ BP = 0.5292$$

문장의 길이가 짧을수록 작은 수를 부여하여 문장 길이에 대한 페널티 적용

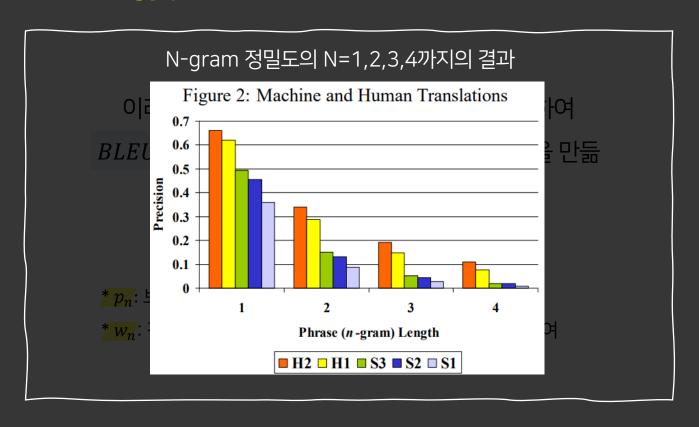
### Brevity Penalty (BP)

이러한 BP와 보정된 N-gram 정밀도를 종합하여 
$$BLEU = BP \times \exp(\sum_{n=1}^{N} w_n \log p_n)$$
 라는 식을 만듦 (일반적으로 N=1,2,3,4)

\*  $p_n$ : 보정된 N-gram 정밀도

\* $w_n$ : 각 N-gram에 대한 가중치, 보통 모두 같은 가중치를 부여

## Brevity Penalty (BP) 을 한 다음 더해주는 이유가 무엇일까?



N이 증가할수록 정밀도는 지수적으로 감소함 이를 보정해주기 위해 log를 취해줌

# 다음 주 예고

1. Attention

2. Transformer

3. Pretrained Model

# 감사합니다