

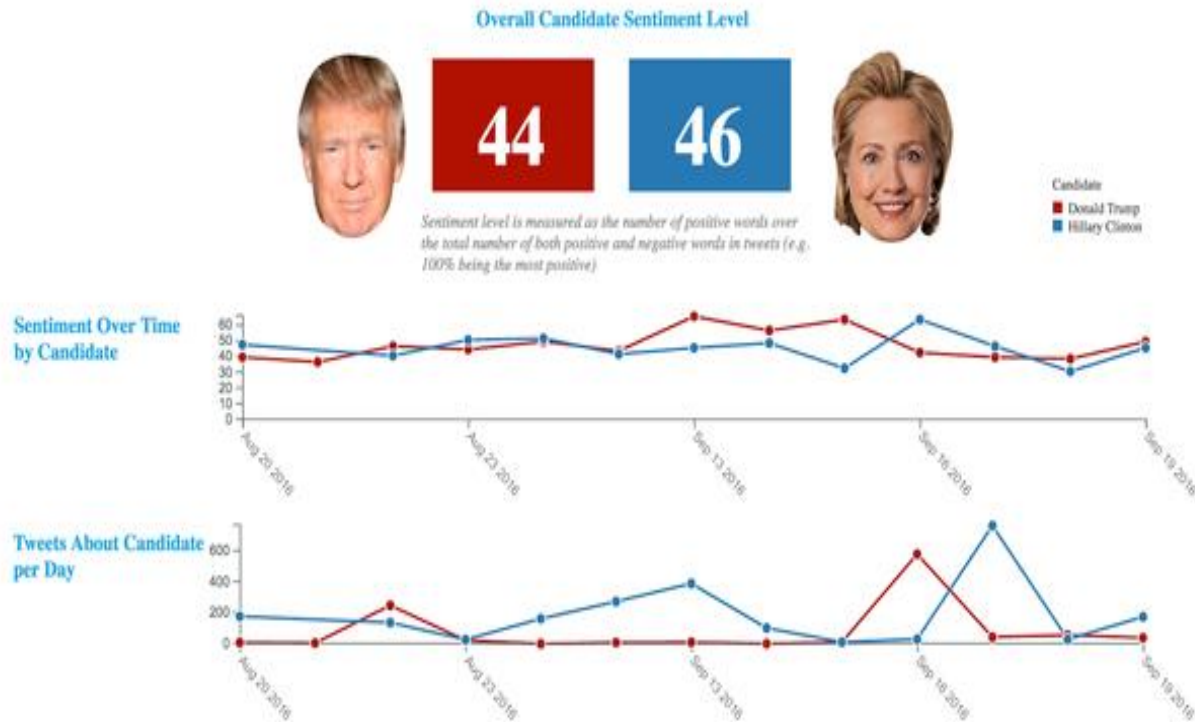
# **PNU** Industrial Data Science

# **Text mining and NLP**

비정형 데이터의 분석

# What is big data?

- Why does everybody want to know about big-data?



*The Trump campaign paid Cambridge Analytica more than \$6 million to help it target voters through ads on Facebook.*

# Contents

산업데이터과학은 산업현장에서 수집된 데이터를 분석하는데 필요한 기초 소양을 강의합니다.

**01**

Unstructured data

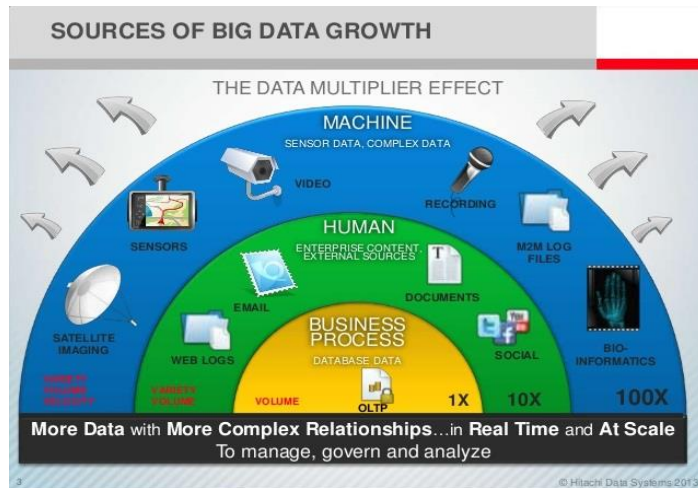
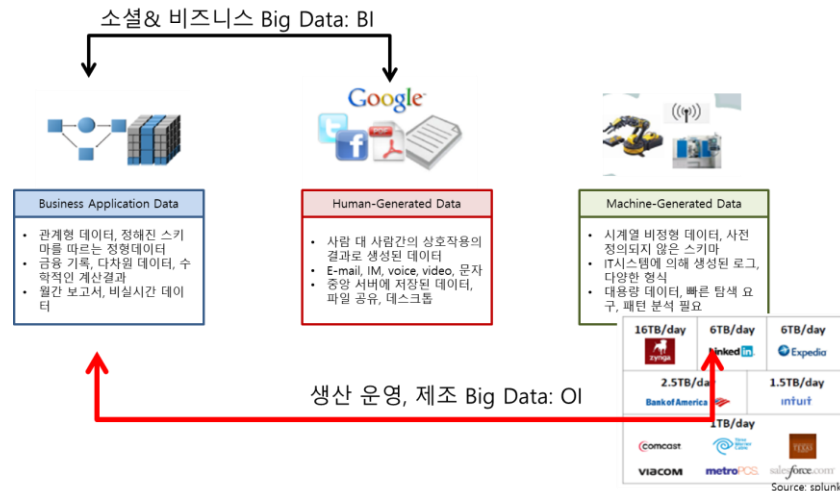
**02**

Text mining

**03**

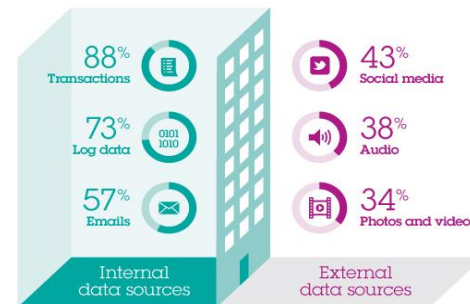
Natural Language  
Processing

# 빅데이터의 종류



## Where does big data come from?

Most big data efforts are currently focused on analyzing internal data to extract insights. Fewer organizations are looking at data outside their firewalls, such as social media.



IBM.

Source: "Capitalize on Big data through Hitachi Innovation", 2013

# 데이터의 종류에 따른 분석



- 정형화되지 않은 데이터
- 미리 정의된 데이터 모델(구조)을 가지고 있지 않은 데이터

예

아주 많은 양의 데이터를 가지고, 구조와 형태가  
다르고 정형화되지 않은 문서, 영상, 음성 등

» 책, 저널, 문서, 메타데이터, 건강 기록, 오디오, 비디오,  
아날로그 데이터, 이미지, 파일, 이메일 메시지,  
웹페이지, 워드프로세스 문서 등

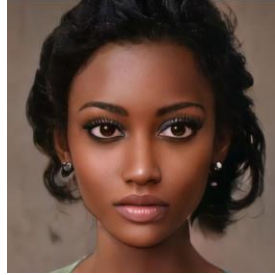
# 비정형 데이터

- 텍스트
  - 텍스트 마이닝, 자연어 처리
- 이미지
  - 이미지 분석
- 음성과 영상
  - 영상 처리
- 로그파일
  - 프로세스 마이닝

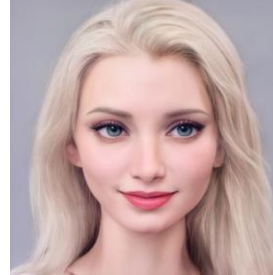
# 다음은 누구의 사진일가요?



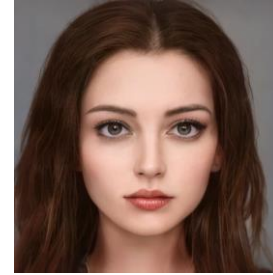
아리엘공주 in 인어공주



티아나공주 in 공주와 개구리



엘사 공주 in 겨울왕국



벨 공주 in 미녀와 야수



이두나 왕비 in 겨울왕국



모아나 in 모아나



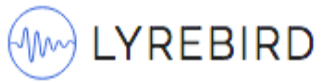
안나 in 겨울왕국





This clip of  
President  
Obama  
talking is fake

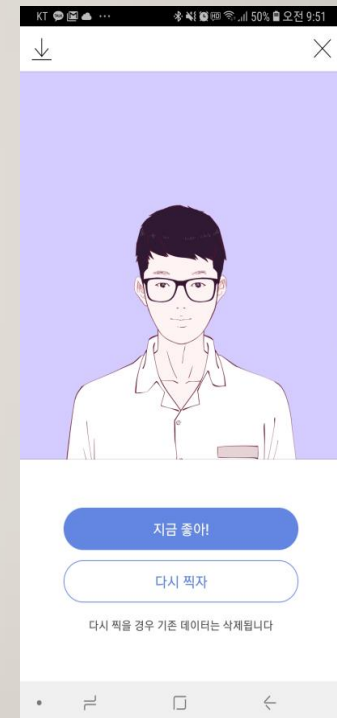




Haeundae Beach is the most famous beach in Busan. The white sand beach is roughly 1.5 kilometer long, over a 30~50 meter wide area, creating a beautiful coastline before a shallow bay, making Haeundae Beach perfect for swimming.



- GNN의 활용 예시
  - N사의 Webtoon에 활용

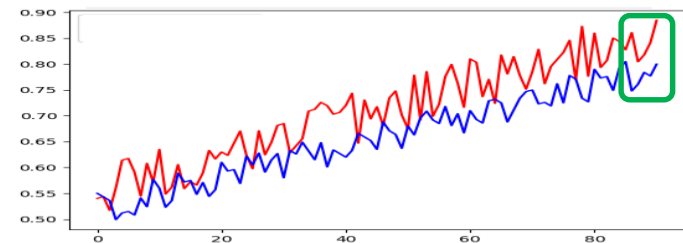


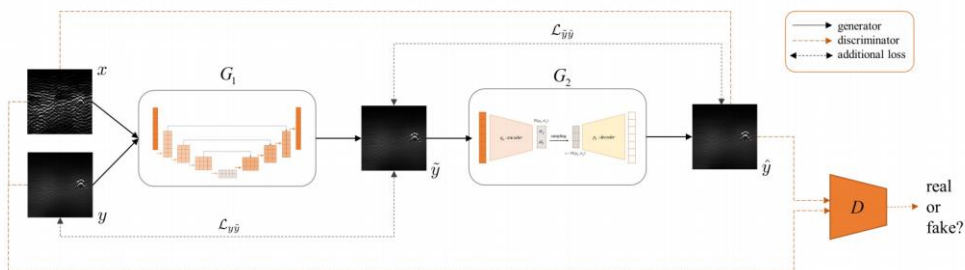
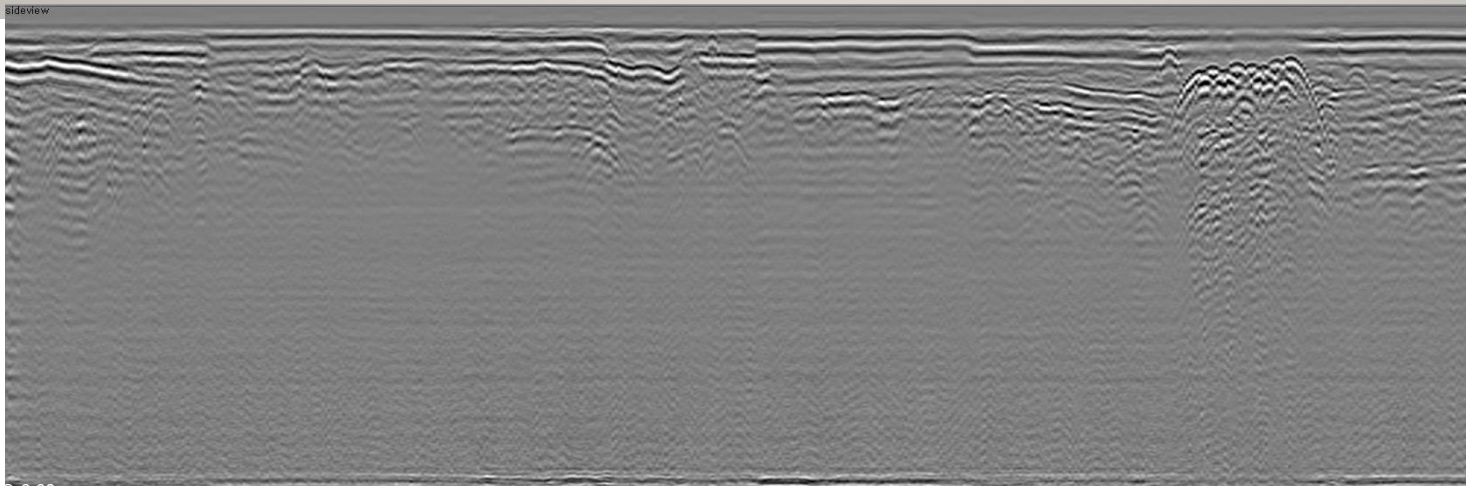


## MBC News '한우 등급 속 여 판 정육점 적발'



4 Layers 4 Classes Training and validation accuracy





Without noise

Model	MSE
DAE	0.0464
DVAE	0.0442
GAFN	0.0048

Noise = 0.5

Model	MSE
DAE	0.0473
DVAE	0.0457
GAFN	0.0155

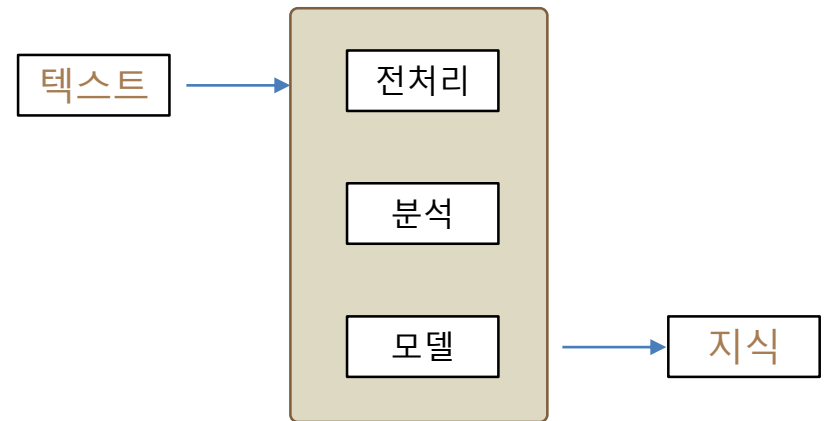
# 비정형 데이터의 마이닝

- 비정형 데이터의 정형화
  - 텍스트 마이닝
    - 데이터 구조로 변화: Meta data의 활용



# 텍스트 마이닝

- 대규모 문서에서
  - 정보추출, 연계성 파악, 분류, 군집화, 요약 등을 시행
  - 텍스트에서 지식을 발견

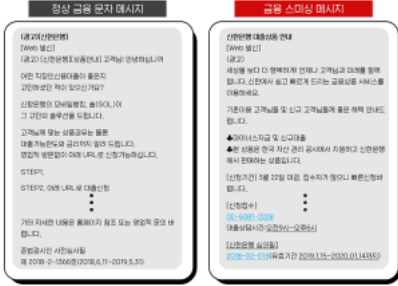


# 사례

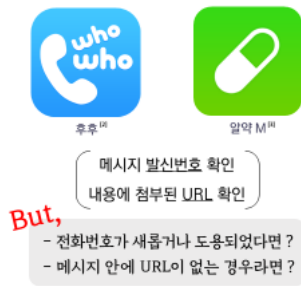
01 연구배경 및 필요성 02 주제 소개 03 수행 내용 04 최종 결과물 05 추진 일정 및 역할분담 06 향후 추진 사항

## 01 연구배경 및 필요성

### 대표적 전자금융사기 - 스미싱, Smishing



### 기존, 스미싱 판별 방식



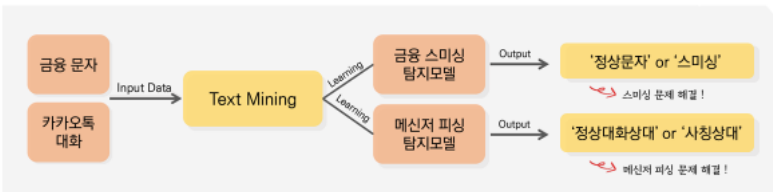
5 [1] 지다영 포스터 [2] WhoWhod&Company

01 연구배경 및 필요성 02 주제 소개 03 수행 내용 04 최종 결과물 05 추진 일정 및 역할분담 06 향후 추진 사항

## 02 주제 소개

Point : 접근유형이 다양하게 변하여도 내용은 변하지 않는다.

“문제에 대한 해결책을 Text 그 자체에서 찾아보자!”



9

01 연구배경 및 필요성 02 주제 소개 03 수행 내용 04 최종 결과물 05 추진 일정 및 역할분담 06 향후 추진 사항

## 01 연구배경 및 필요성



### 새로운 전자금융사기 - 메신저 피싱

- 가족, 주변 지인의 프로필 사진과 이름을 도용
- 전자기기 사용이 익숙한 중장년층을 타겟 설정
- 다급한 상황 연출 + 타인 계좌로의 송금 요청
- 반복적인 보이스트록 연결을 통한 외부로의 연락 방해

“해결방법은 단지 스스로 조심하기”

6

01 연구배경 및 필요성 02 주제 소개 03 수행 내용 04 최종 결과물 05 추진 일정 및 역할분담 06 향후 추진 사항

## 04 최종 결과물 - 모델 성능 평가, 메신저 피싱

3개의 모델 모두, 약 83%의 특이도와 약 76%의 정확도의 성능을 보여준다. 이것은 어느 수준 이상이 되면 모델의 선택이 아닌 데이터 전처리를 재고민해야 함을 알 수 있다.



25

# 웹마이닝

- 웹 콘텐츠 분석
  - 웹페이지내의 콘텐츠로 부터 데이터, 정보, 지식을 추출함
- 웹 행위 분석
  - 웹사이트의 페이지간의 연결구조를 분석



# 사례

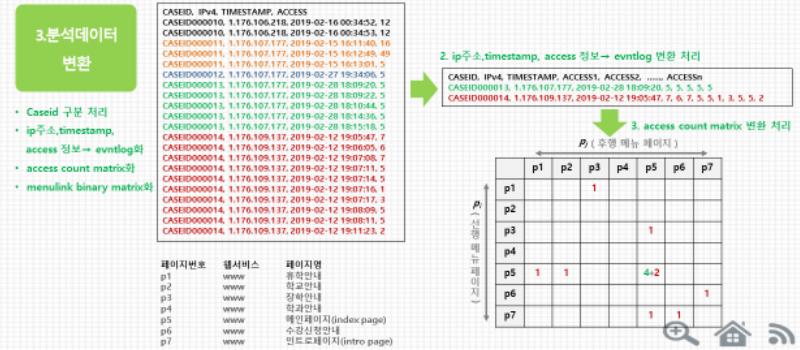
## 2. 연구방법

### 1) 연구 분석 절차



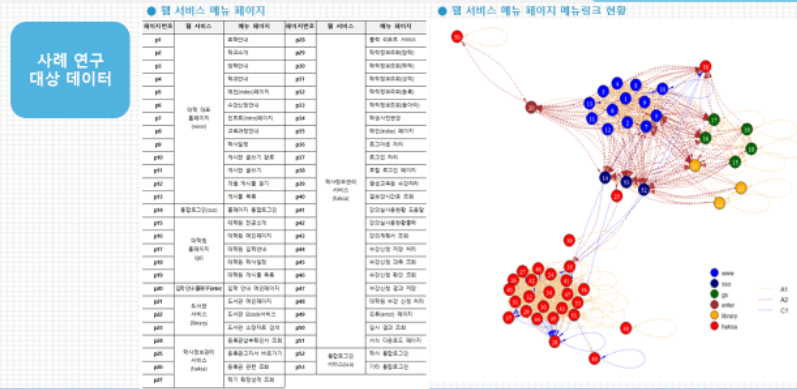
## 2. 연구방법

### 2) 연구 분석 처리



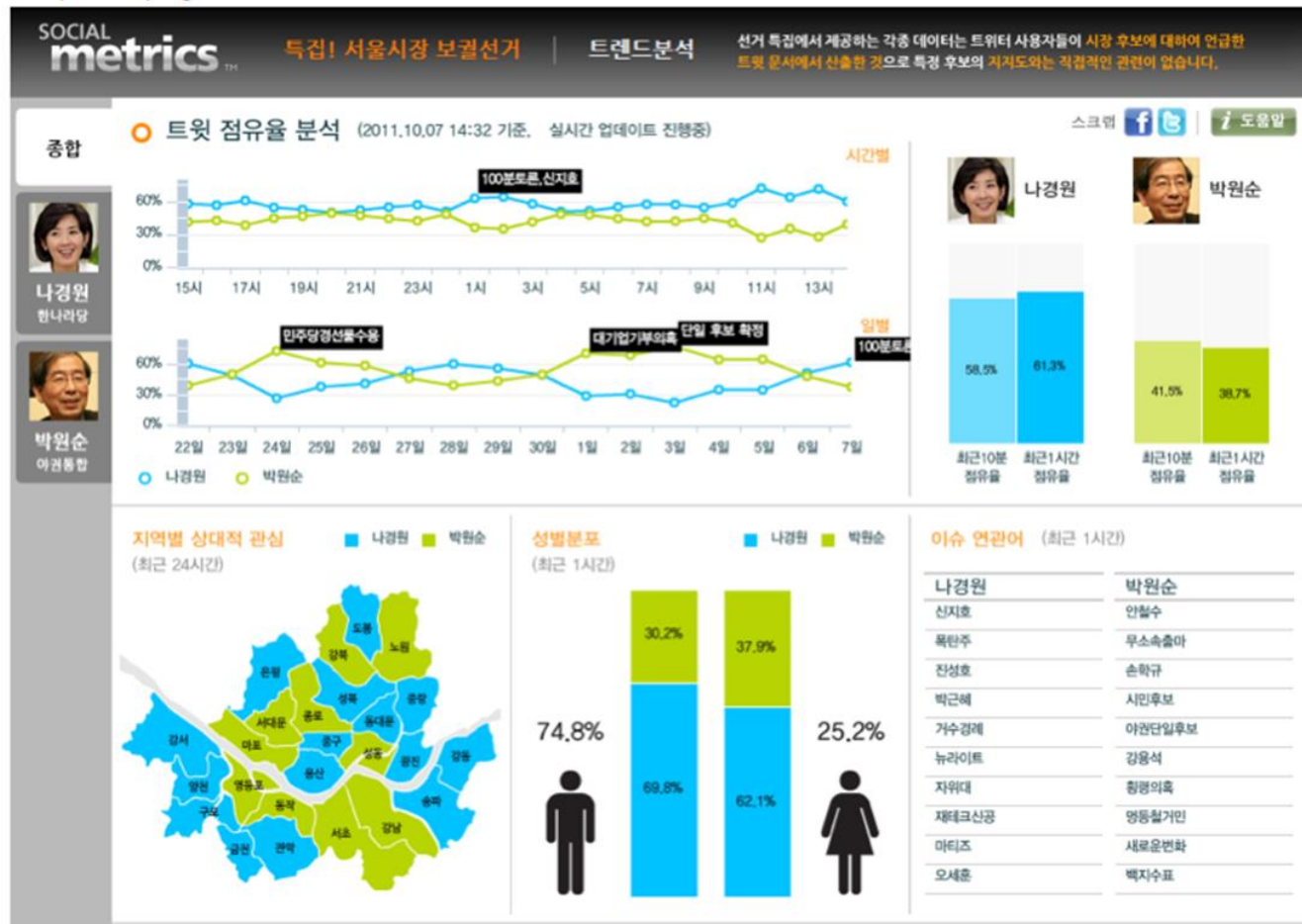
## 3. 사례연구

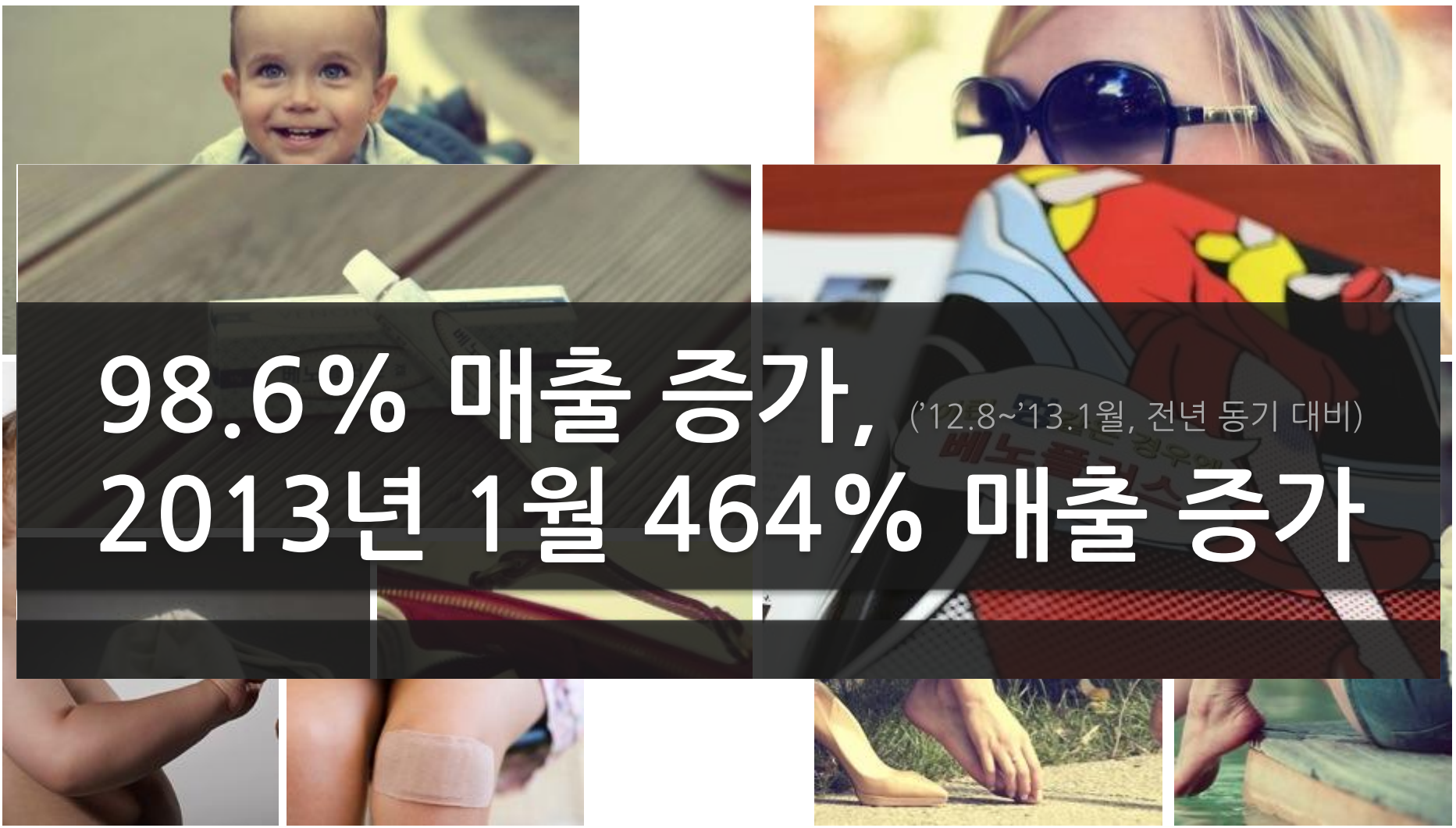
### 1) 연구 데이터



# Opinion Mining

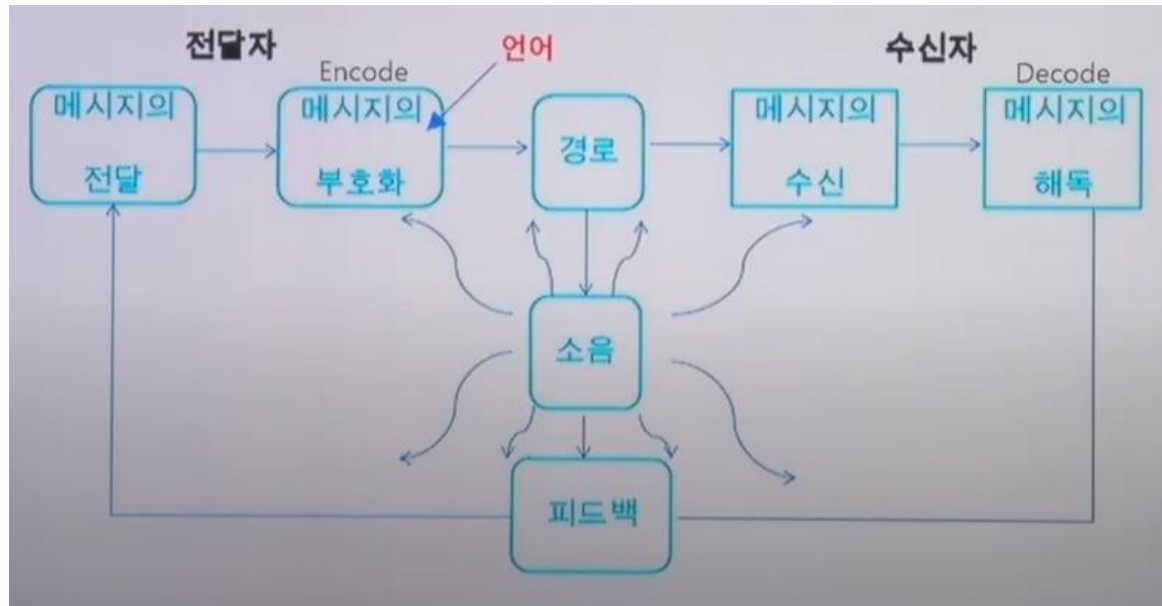
<http://campaign.socialmetrics.co.kr/>





98.6% 매출 증가, (’12.8~’13.1월, 전년 동기 대비)  
2013년 1월 464% 매출 증가

# NLP



출처:[토큰ON세미나] 자연어 언어모델 'BERT' 1강 - 자연어 처리 (NLP) | T아카데미  
<https://www.youtube.com/watch?v=qlxrXX5uBoU&list=PLCTJCvcltJv0mSXDydtOLB2InrtksJEQ&index=26>

# NLP Problem

- Doc. Classification
- Grammar correction
- Information Extraction
- Voice recognition
- Information search
- Abstracting
- Translation
- Q&A
- Machine interpretation
- Chatbot
- Morpheme analysis
- Sentiment analysis
- Intention analysis

**01** | 자연어 처리 (Natural Language Processing, NLP)

다양한 자연어 처리 Applications

• 문서 분류	• 형태소 분석
• 문법, 오타 교정	• 개체명 분석
• 정보 추출	• 구문 분석
• 음성 인식결과 보정	• 감성 분석
• 음성 합성 텍스트 보정	• 관계 추출
• 정보 검색	• 의도 파악
• 요약문 생성	
• 기계 번역	
• 질의 응답	
• 기계 독해	
• 챗봇	

출처:[토크ON세미나] 자연어 언어모델 'BERT' 1강 - 자연어 처리 (NLP) | T아카데미  
<https://www.youtube.com/watch?v=qlxrXX5uBoU&list=PLCTJCvcltJv0mSXDydtOLB2InrtnksJEQ&index=26>



# Classification and NLP

- Many NLP problems belong to classification

Question?  
Request?  
Rejection?  
Approval?

Predicate?  
Adjective?  
Noun?

Positive?  
Neutral?  
Negative?

Noun?  
Adverb?  
Proper noun?  
Verb?

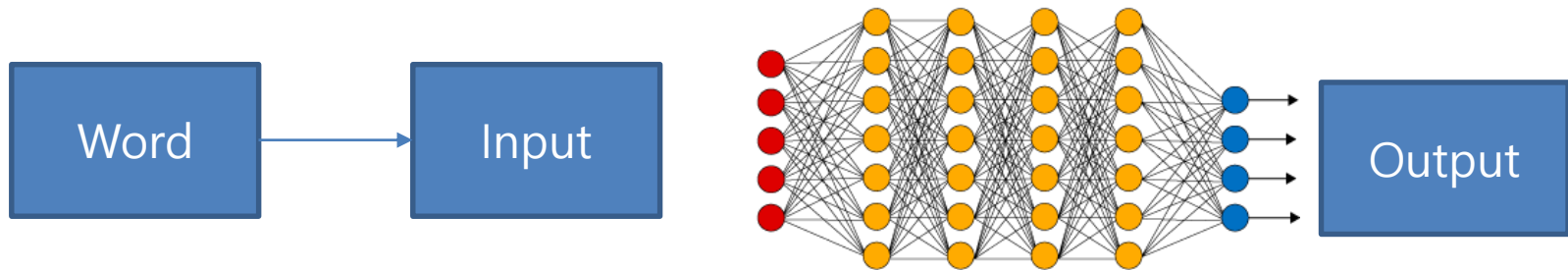
Organization?  
Person?  
Position?



출처:[토크ON세미나] 자연어 언어모델 'BERT' 1강 - 자연어 처리 (NLP) | T아카데미  
<https://www.youtube.com/watch?v=qlxrXX5uBoU&list=PLCTJCvcltJv0mSXDydtOLB2lnrtksJEQ&index=26>

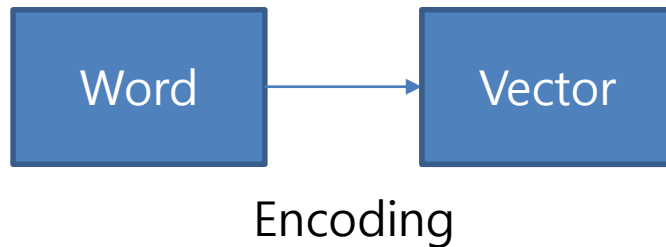
# In order to analyze text...

- Can words be input to deep learning?
- We need to convert words to input of deep learning



# What is word2Vec?

- Converting words to input of deep learning
- Encoding
  - Converting text to number
  - Number is represented as vector



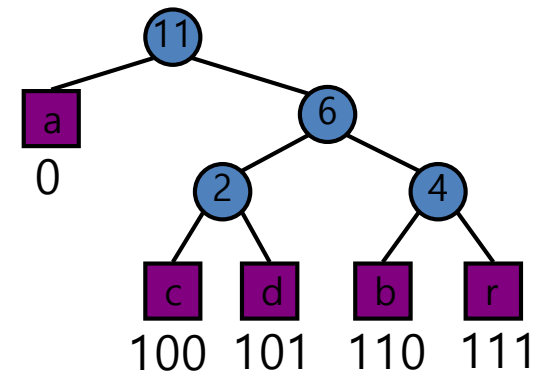


# Encoding

- How to convert text into number
  - "Hyerim loves Wonkyung"
  - Hyerim: 0, loves: 1, Wonkyung: 2
- Frequency based encoding: Huffman Encoding
- Label encoding
  - Hyerim: 1, loves: 2, Wonkyung: 3...
- One-hot-ecoding
  - Hyerim: [1, 0, 0], loves: [0, 1, 0], Wonkyung: [0, 0, 1]

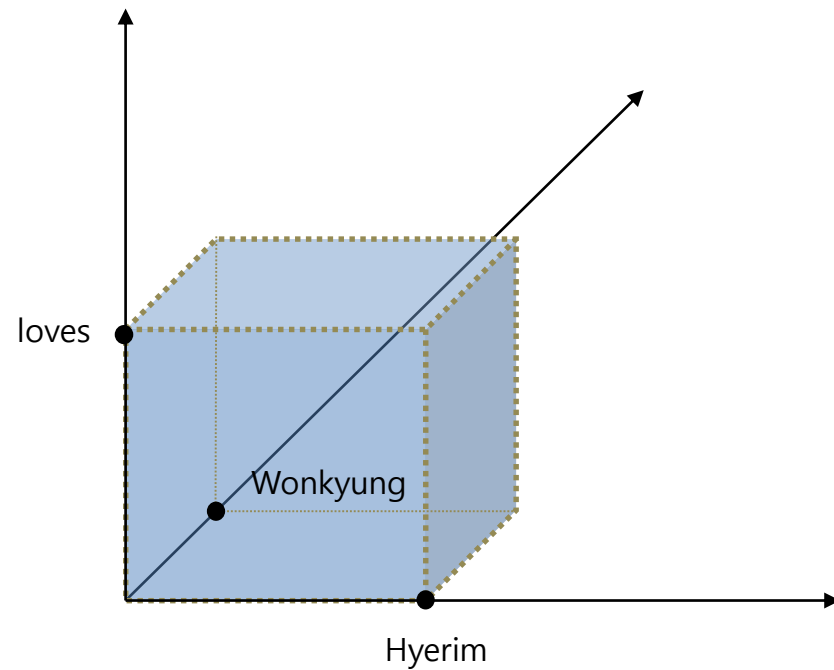
$X = \text{abracadabra}$

011011101000101...



# Problems of OHE

- No similarity
  - Cannot represent similarity b/w words
  - Cosine similarity is always 0

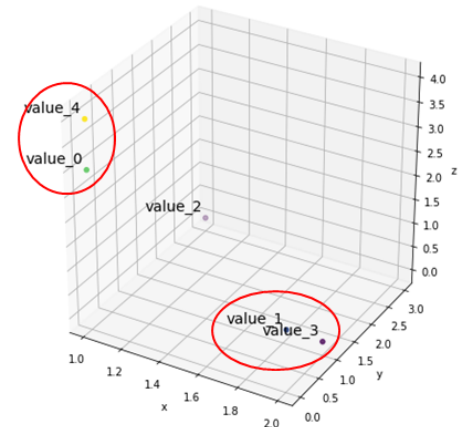


# Embedding

- Embedding is dense vector (One category is not one dimension)
  - With similarity

Value of Categorical Feature	Embedding
value_0	[1.0, 0.0, 3.0]
value_1	[2.0, 0.0, 1.0]
value_2	[1.0, 3.0, 0.0]
value_3	[2.0, 1.0, 0.0]
value_4	[1.0, 0.0, 4.0]

$$p(w_c|w) = \frac{\exp(\mathbf{w} \cdot \mathbf{w}_c)}{\sum_i \exp(\mathbf{w} \cdot \mathbf{w}_i)}$$



# Word2Vec is word embedding

- Similarity comes from neighboring relations
  - “king brave man”, “queen beautiful woman”

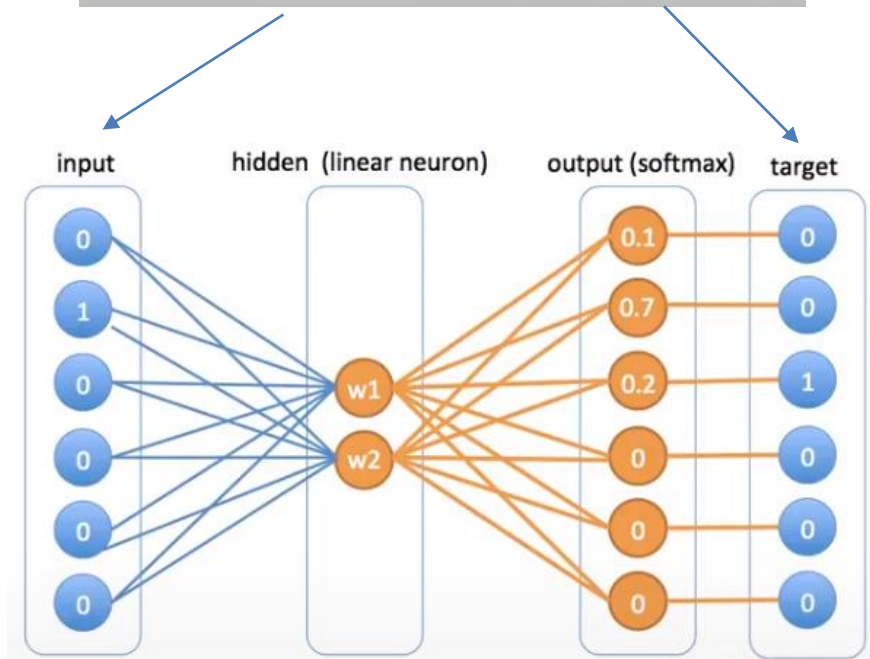
word	neighbor
king	brave
brave	king
brave	man
man	brave
queen	beautiful
beautiful	queen
beautiful	woman
woman	beautiful

Windows size =1

# Word2Vec using NN

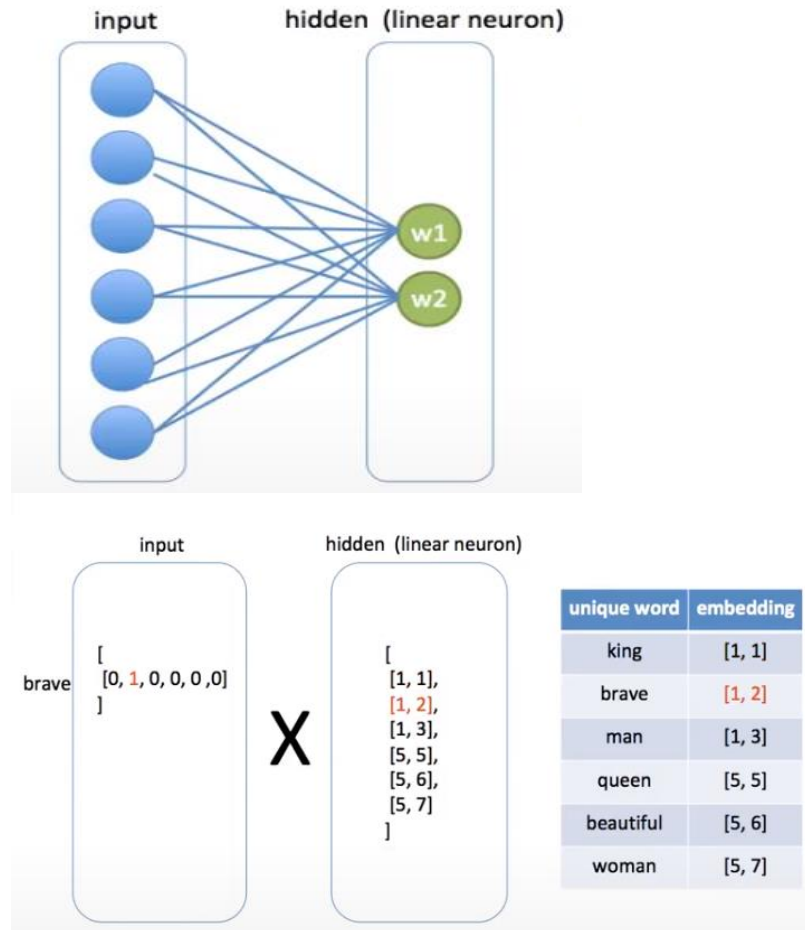
- Training data (Time window = 2)

word	word one hot encoding	neighbor	neighbor one hot encoding
king	[1, 0, 0, 0, 0, 0]	brave	[0, 1, 0, 0, 0, 0]
king	[1, 0, 0, 0, 0, 0]	man	[0, 0, 1, 0, 0, 0]
brave	[0, 1, 0, 0, 0, 0]	king	[1, 0, 0, 0, 0, 0]
brave	[0, 1, 0, 0, 0, 0]	man	[0, 0, 1, 0, 0, 0]
man	[0, 0, 1, 0, 0, 0]	king	[1, 0, 0, 0, 0, 0]
man	[0, 0, 1, 0, 0, 0]	brave	[0, 1, 0, 0, 0, 0]
queen	[0, 0, 0, 1, 0, 0]	beautiful	[0, 0, 0, 0, 1, 0]
queen	[0, 0, 0, 1, 0, 0]	woman	[0, 0, 0, 0, 0, 1]
beautiful	[0, 0, 0, 0, 1, 0]	queen	[0, 0, 0, 1, 0, 0]
beautiful	[0, 0, 0, 0, 1, 0]	woman	[0, 0, 0, 0, 0, 1]
woman	[0, 0, 0, 0, 0, 1]	queen	[0, 0, 0, 1, 0, 0]
woman	[0, 0, 0, 0, 0, 1]	beautiful	[0, 0, 0, 0, 1, 0]



# Word2Vec

- Word2Vec is hidden layer

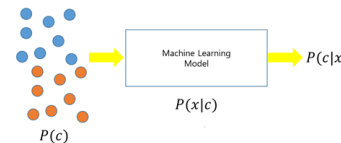


[http://github.com/minsuk-heo/python\\_tutorial/blob/master/data\\_science/nlp/word2vec\\_tensorflow.ipynb](http://github.com/minsuk-heo/python_tutorial/blob/master/data_science/nlp/word2vec_tensorflow.ipynb)

# Bayesian NLP

- Recall ML classification and Bayesian
- From text we want to classify
  - Is the mail spam or not?
  - Use Naïve Bayes
- Problem of NB
  - Cannot consider the order in the text

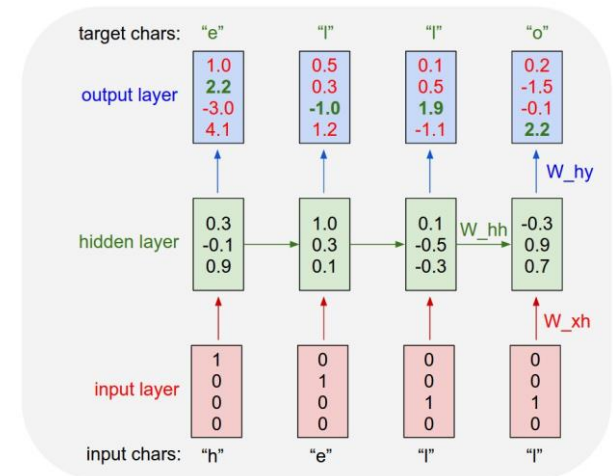
- Given data, decide a class
- Machine Learning은 data set이 주어졌을 때, 특정 사건 or 가설의 확률을 높여줄 수 있는 최적의 모델을 찾는 것을 목적으로 함
  - 베이지안 모델이 주어진 정보를 update 해 가며 최적의 사후확률을 계산하는 방식과 일맥상통



$$P_{nb}(C_1 | \dots, x_p)$$
$$= \frac{P(C_1)[P(x_1|C_1)P(x_2|C_1)\cdots P(x_p|C_1)]}{P(C_1)[P(x_1|C_1)P(x_2|C_1)\cdots P(x_p|C_1)] + \cdots + P(C_m)[P(x_1|C_m)P(x_2|C_m)\cdots P(x_p|C_m)]}$$

# NLP and time series

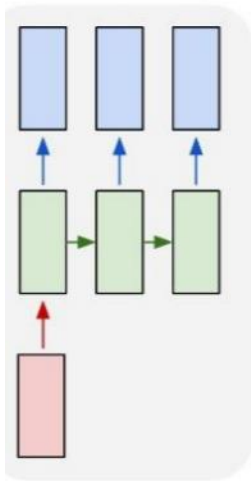
- Previous character influences next character



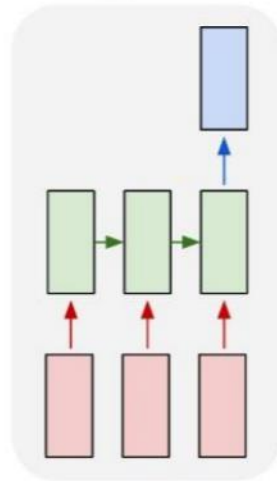


- RNN can solve various types of time series problems

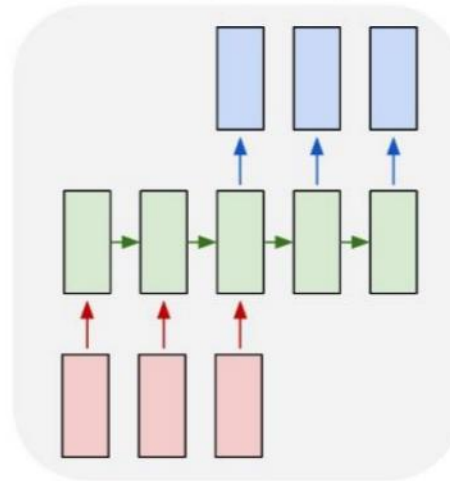
one to many



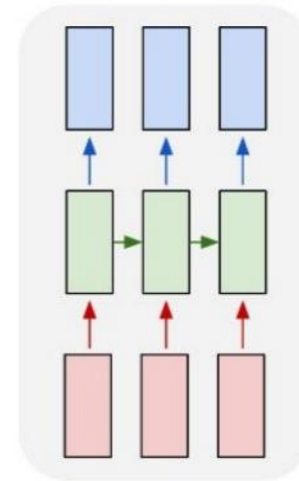
many to one



many to many

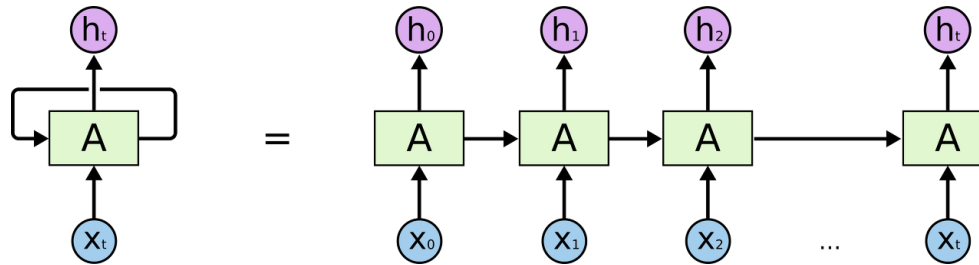


many to many



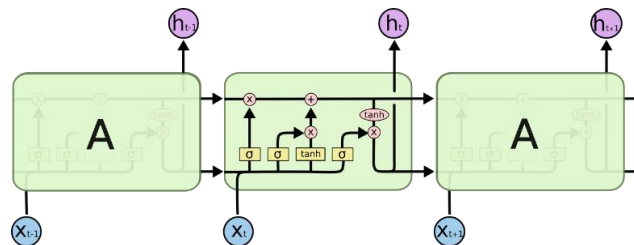
# RNN and LSTM

- RNN



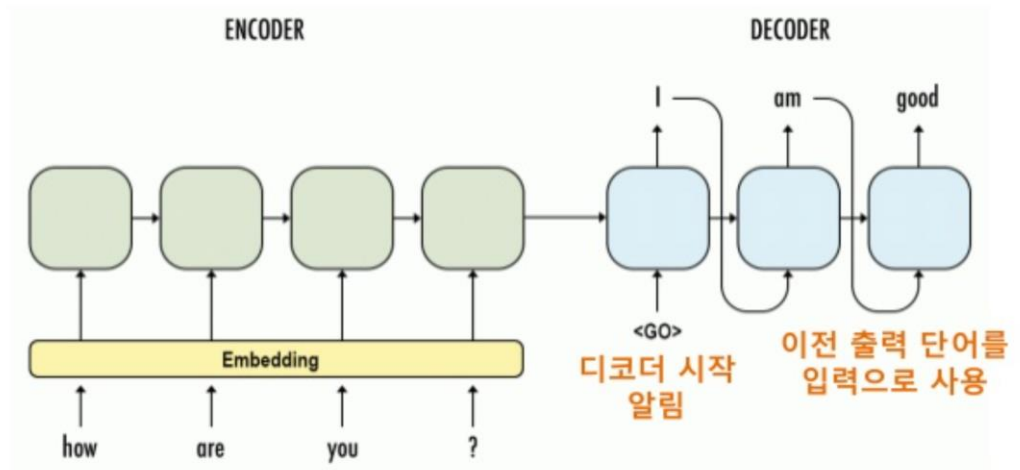
- LSTM

"the clouds are in the *sky*"  
"I grew up in France... I speak fluent *French*"



# Sequence to Sequence

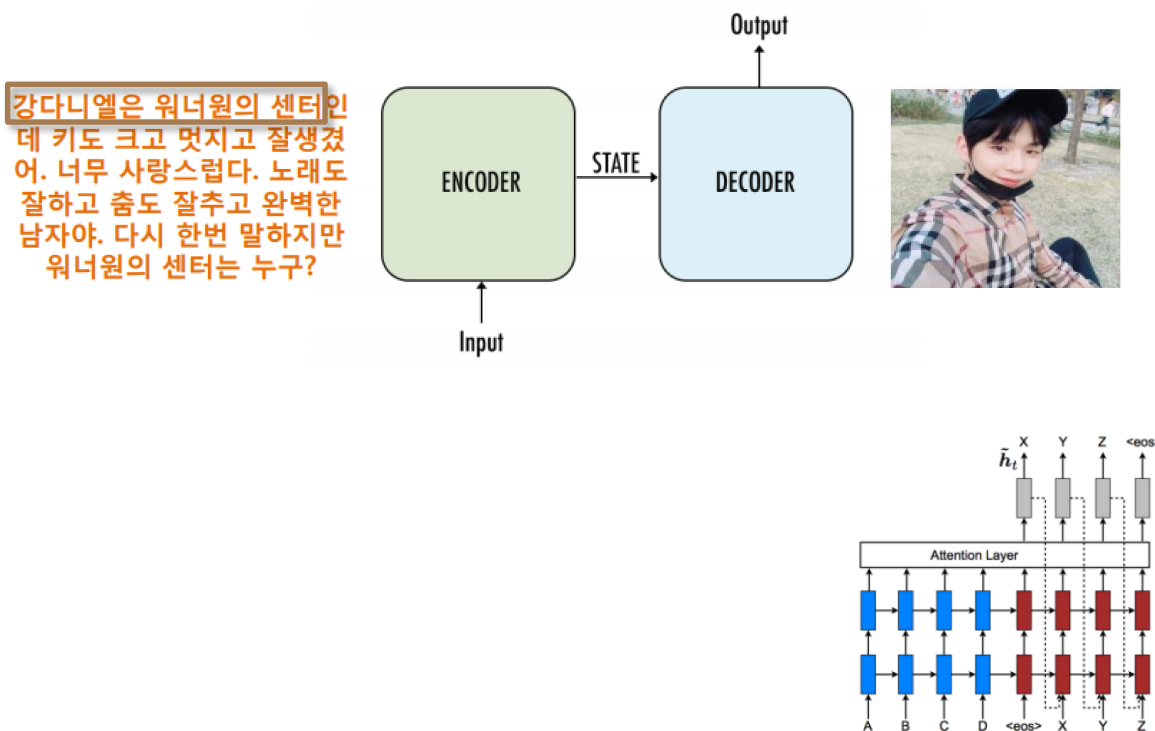
- RNN base approach
  - Cannot process whole text input



출처: 우종하, "딥러닝 자연어처리:RNN에서 BERT까지"  
[https://www.slideshare.net/deepseaswjh/rnn-bert?from\\_action=save](https://www.slideshare.net/deepseaswjh/rnn-bert?from_action=save)

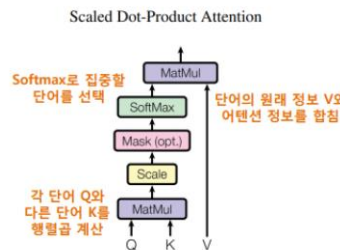
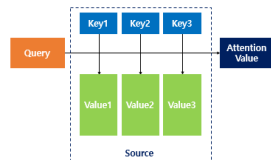
# Attention

- Seq2Seq
  - In Encoder, there is loss of information

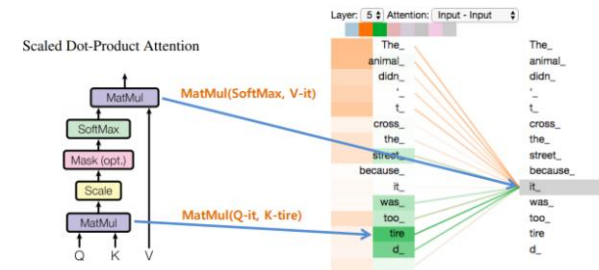


# Transformer

- By google, 2017
- LSTM is not needed
- Self attention



Input	Thinking	Machines
Embedding	$x_1$	$x_2$
Queries	$q_1$	$q_2$
Keys	$k_1$	$k_2$
Values	$v_1$	$v_2$
Score	$q_1 \cdot k_1 = 112$	$q_1 \cdot k_2 = 96$
Divide by $8 (\sqrt{d_k})$	14	12
Softmax	0.88	0.12



## Attention Is All You Need

Ashish Vaswani<sup>\*</sup>  
Google Brain  
avaswani@google.com

Noam Shazeer<sup>\*</sup>  
Google Brain  
noam@google.com

Niki Parmar<sup>\*</sup>  
Google Research  
nikip@google.com

Jakob Uszkoreit<sup>\*</sup>  
Google Research  
uszk@google.com

Llion Jones<sup>\*</sup>  
Google Research  
llion@google.com

Aidan N. Gomez<sup>†</sup>  
University of Toronto  
aidan@ca.toronto.edu

Lukasz Kaiser<sup>\*</sup>  
Google Brain  
lukaszkaiser@google.com

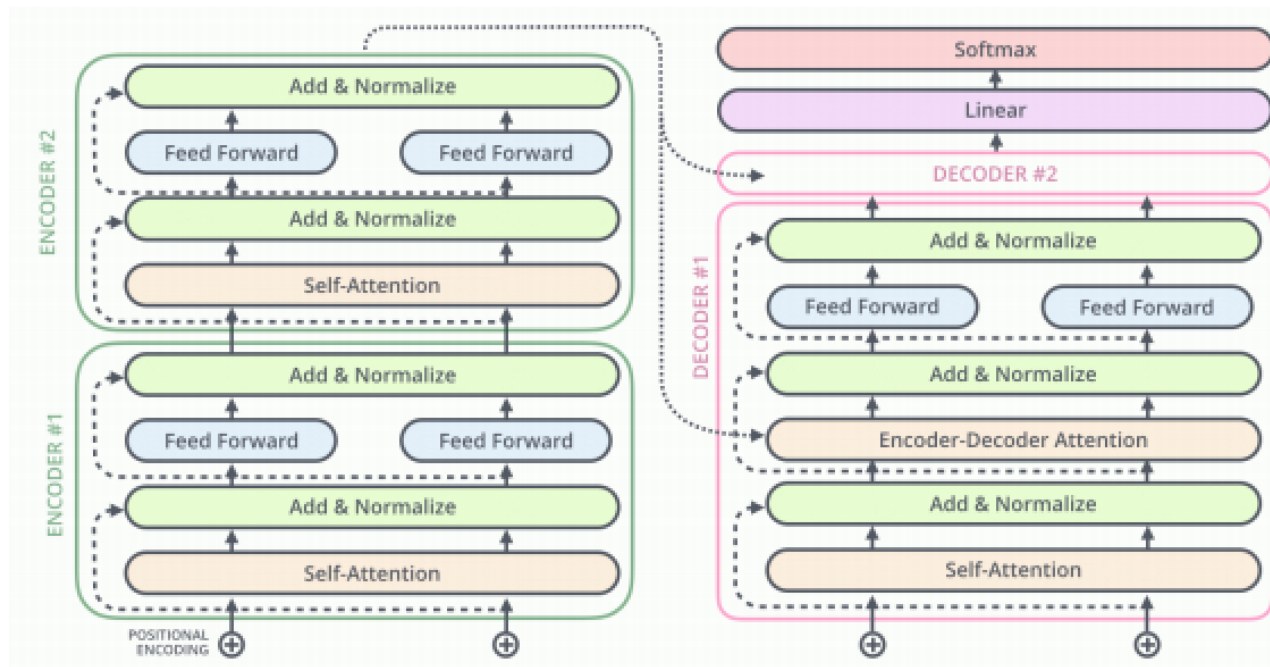
Illia Polosukhin<sup>\*</sup><sup>‡</sup>  
illia.polosukhin@gmail.com

### Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.8 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature. We show that the Transformer generalizes well to other tasks by applying it successfully to English constituency parsing both with large and limited training data.

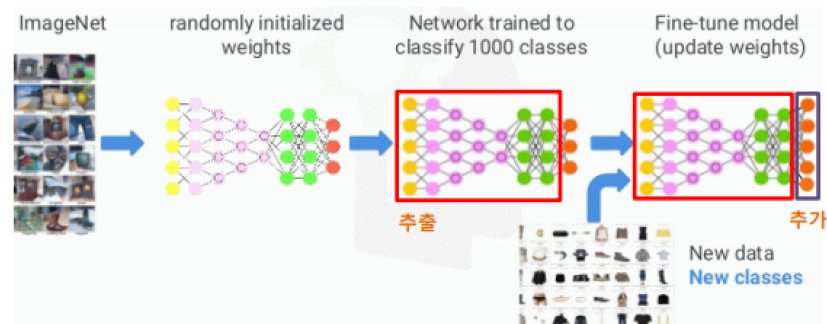
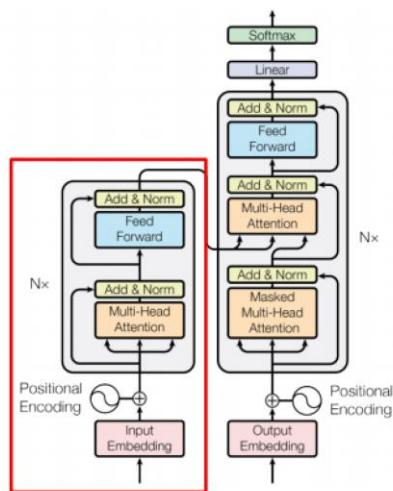
# Structure of 'Transformer'

- Nested structure of encoder and decoder



# BERT

- By google, 2018
- Pretrained model for transfer learning
- Using transformer
  - BERT Base: 12 transformers
  - BERT Large: 24 transformers



## SQuAD competition ranking

Rank	Model	EM	F1
	Human Performance Stanford University (Rajpurkar & Jia et al. '18)	86.831	89.452
1	BERT finetune baseline (ensemble) Anonymous	83.536	86.096
2	Lunet + Verifier + BERT (ensemble) Layer 6 AI NLP Team	83.469	86.043
3	Lunet + Verifier + BERT (single model) Layer 6 AI NLP Team	82.995	86.035
4	PAML+BERT (single model) PINGAN GammaLab	82.577	85.603
5	AoA + DA + BERT (ensemble) Joint Laboratory of HIT and iFLYTEK Research	82.374	85.310