

Deep Learning for NLP

- Word Embedding -

Ko, Youngjoong

September 9, 2015

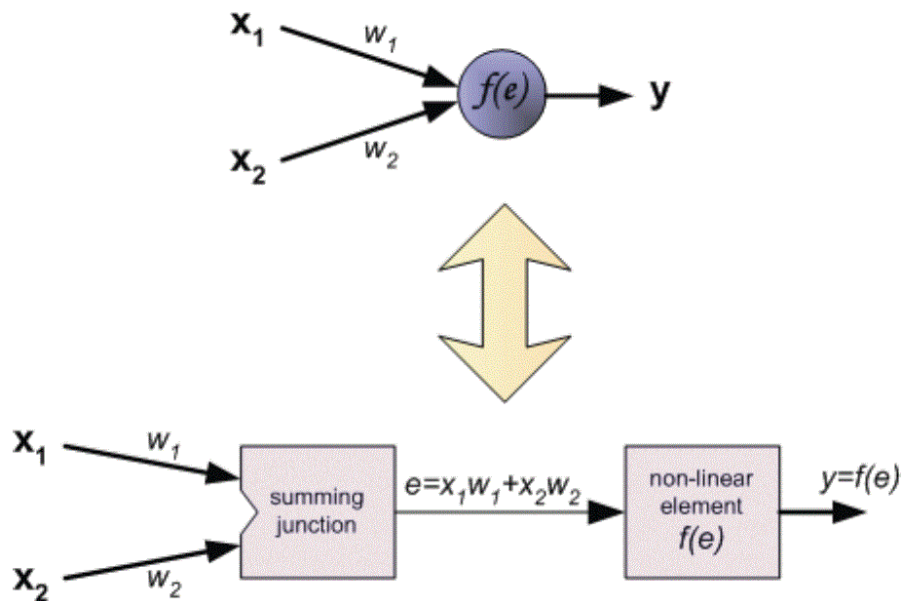
Dept. of Computer Engineering,
Dong-A University

Contents

- 1. Basic Concepts of Neural Network (NN)**
- 2. Why do we need Deep Learning?**
- 3. Learning Representation for NLP**
- 4. Tools for Word Embedding**
 - Word2Vector
 - Ranking-based

Basic Concepts of NN

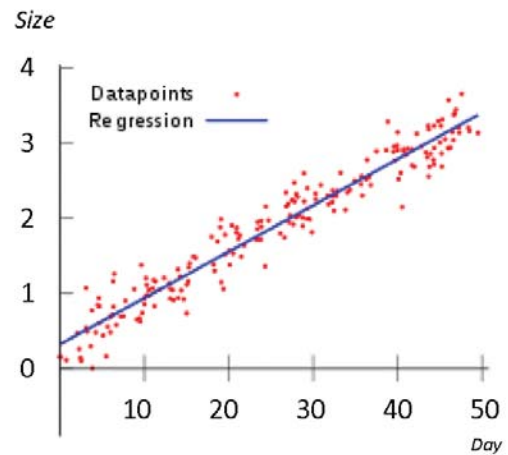
❖ Perceptron



3

Basic Concepts of NN

❖ Illustration Example (Apple Tree)

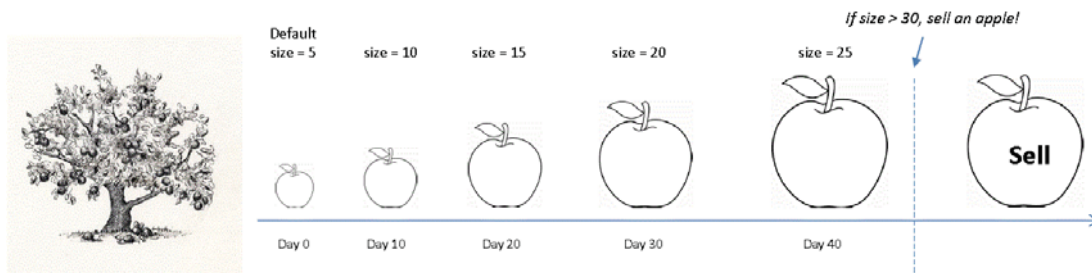


- "어떤 사과나무에 대해서 몇 년에 걸쳐 날짜 별로 사과들의 크기를 측정, 기록"
- 농부는 특정 크기가 넘을 때만 시장에 사과를 내다 팔 수 있다고 할 때,
- Q: 올해 Day -50 에 사과를 내다 팔 수 있을까? 없을까?

4

Basic Concepts of NN

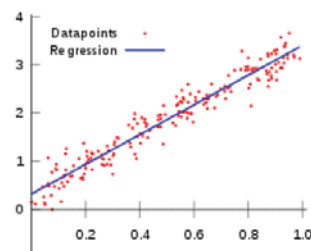
❖ Illustration Example (Apple Tree)



상황 1 : 작년까지 이 사과나무는 위의 경향대로 사과 열매를 맺었다.
조건 : 사과的大小가 30이 넘으면 팔 수 있다.

Question : 올해 Day-50 에 사과를 팔 수 있을까?

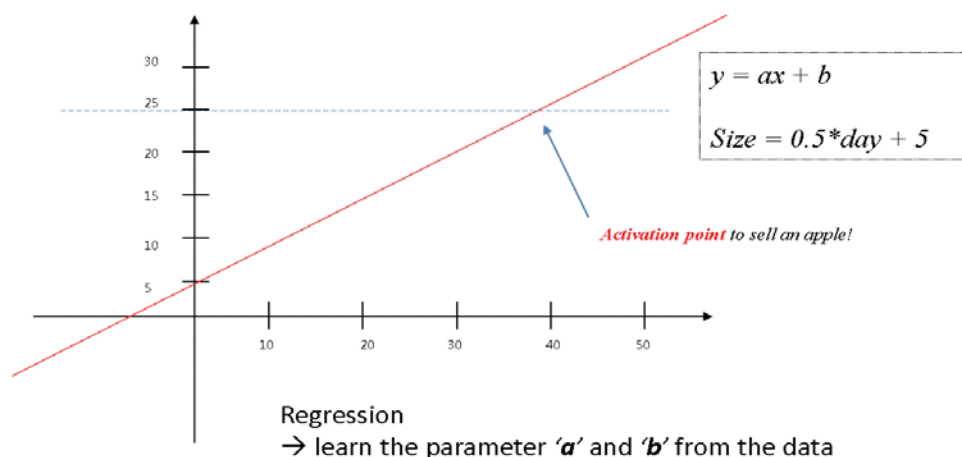
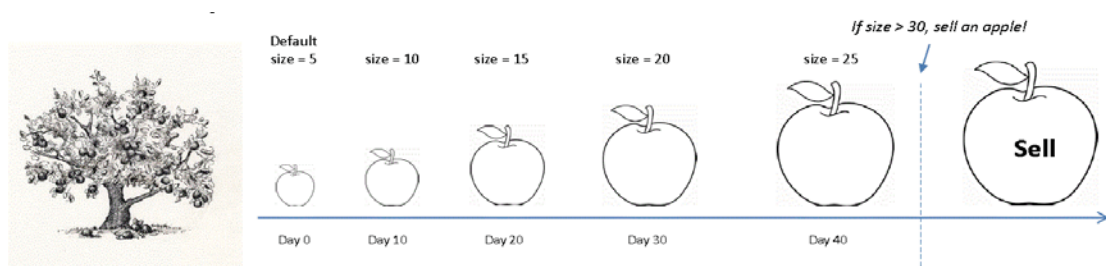
Very Typical **Regression Problem**



5

Basic Concepts of NN

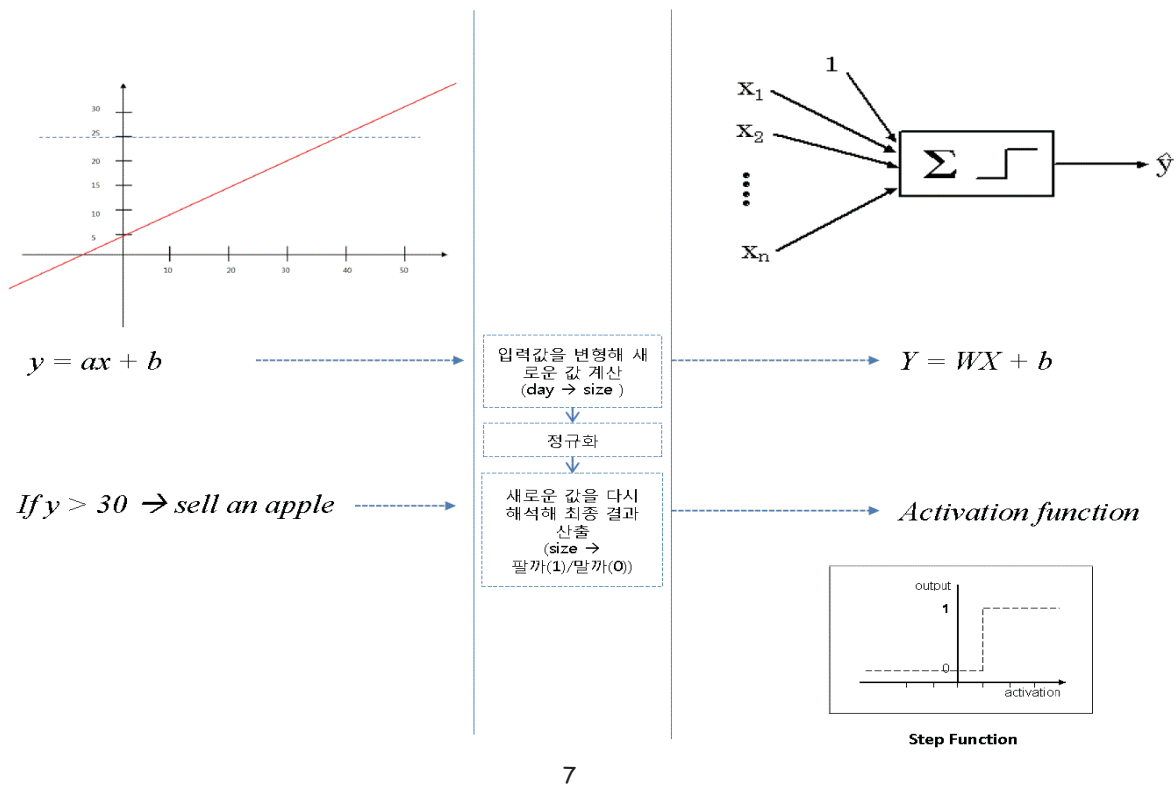
❖ Illustration Example (Apple Tree)



6

Basic Concepts of NN

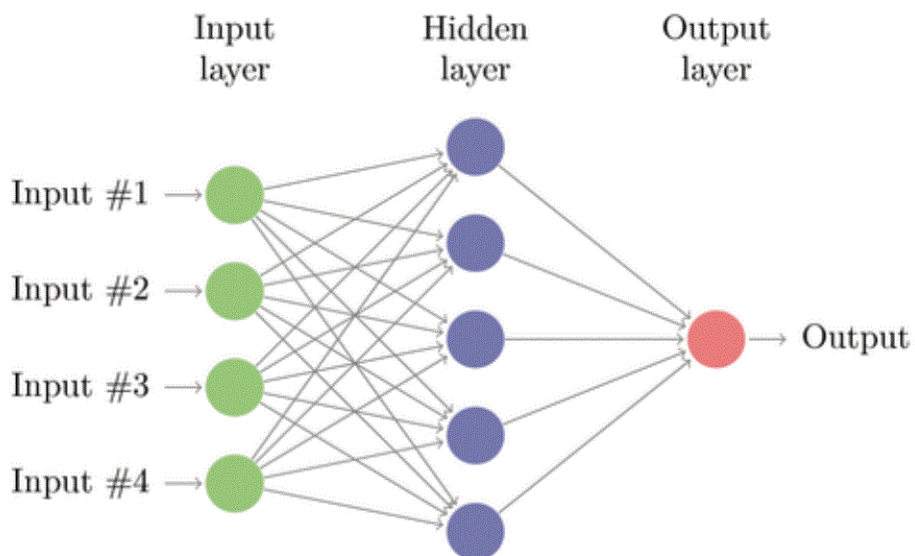
❖ Illustration Example (Apple Tree)



7

Basic Concepts of NN

❖ Multilayer Neural Network



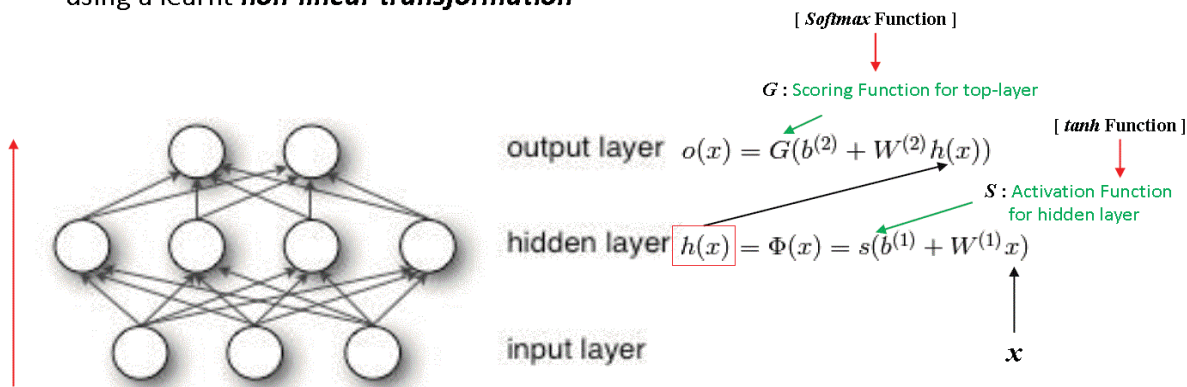
8

Basic Concepts of NN

❖ Multilayer Neural Network

The single-hidden layer Multi-Layer Perceptron (MLP).

An **MLP** can be viewed as a **logistic regressor**, where the input is first transformed using a learnt **non-linear transformation**



$$f : R^D \rightarrow R^L$$

$$f(x) = G(b^{(2)} + W^{(2)}(s(b^{(1)} + W^{(1)}x))),$$

D is the size of input vector x

L is the size of output vector $f(x)$

Feed Forward Propagation

Basic Concepts of NN

❖ Training (Weight Optimization)

$$\theta = \{W^{(2)}, b^{(2)}, W^{(1)}, b^{(1)}\}$$

- How to learn the weights??

“Backpropagation Algorithm”

최종 결과물을 얻고	Feed Forward and Prediction
그 결과물과 우리가 원하는 결과물과의 차이점을 찾은 후	Cost Function
그 차이가 무엇으로 인해 생기는 지	Differentiation (미분)
역으로 내려가면서 추정하여	Back Propagation
새로운 Parameter 값을 배움	Weight Update

Basic Concepts of NN

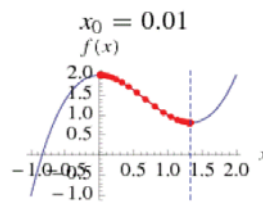
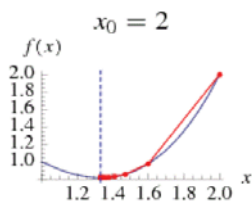
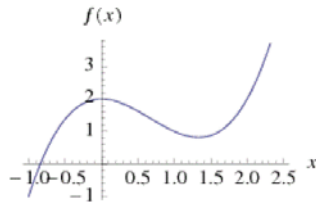
❖ Training (Weight Optimization)

Backpropagation = Backpropagation of errors



Gradient descent procedures are generally used where we want to maximize or minimize n-dimensional functions.

The **gradient** is a vector \mathbf{g} that is defined for any **differentiable** point of a function, that points from this point exactly towards the **steepest ascent** and indicates the gradient in this direction by means of its norm $|\mathbf{g}|$.



$$f(x) = x^3 - 2x^2 + 2$$

$$x_i = x_{i-1} - \epsilon f'(x_{i-1})$$

x_i 가 변화가 없을 때까지 위 수식을 반복

→ 결국 gradient 가 가리키는 방향으로 계속해서 parameter 변화 됨

→ Local Minimum 에 빠질 수 있음

Basic Concepts of NN

❖ Training (Activation Functions)

$\text{sigmoid}(a) = 1/(1 + e^{-a})$... also called '**logistic function**', '**Fermi function**'

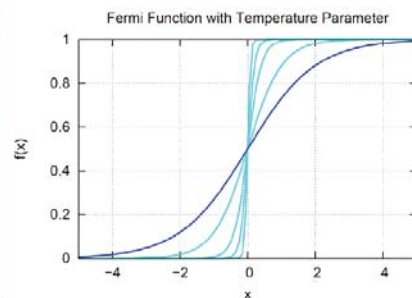
$$f(x) = \frac{1}{1 + e^{-x}}$$

$$\frac{d}{dx}f(x) = f(x)(1 - f(x))$$

$$1 - f(x) = f(-x).$$

$$2f(x) = 1 + \tanh\left(\frac{x}{2}\right).$$

Always positive



$$\text{tanh}(a) = (e^a - e^{-a}) / (e^a + e^{-a})$$

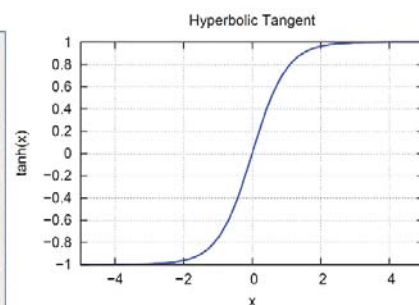
$$f(x) = \frac{\sinh x}{\cosh x} = \frac{e^x - e^{-x}}{e^x + e^{-x}} = \frac{e^{2x} - 1}{e^{2x} + 1}$$

$$e^x = \cosh x + \sinh x$$

and

$$e^{-x} = \cosh x - \sinh x$$

Output = [-1, 1]
Faster Backpropagation



Basic Concepts of NN

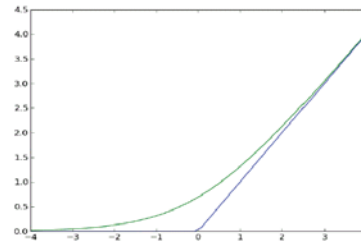
❖ Training (Activation Functions)

Rectified Linear Unit $f(x) = \max(0, x)$

Smooth approximation – “softplus” function

$$f(x) = \log(1 + e^x)$$

$$f'(x) = e^x / (e^x + 1) = 1 / (1 + e^{-x})$$



❖ Scoring Functions (Softmax)

$$\text{softmax}_{ij}(x) = \frac{\exp x_{ij}}{\sum_k \exp(x_{ik})}$$

$$P(Y = i | x, W, b) = \text{softmax}_i(Wx + b)$$

$$= \frac{e^{W_i x + b_i}}{\sum_j e^{W_j x + b_j}}$$

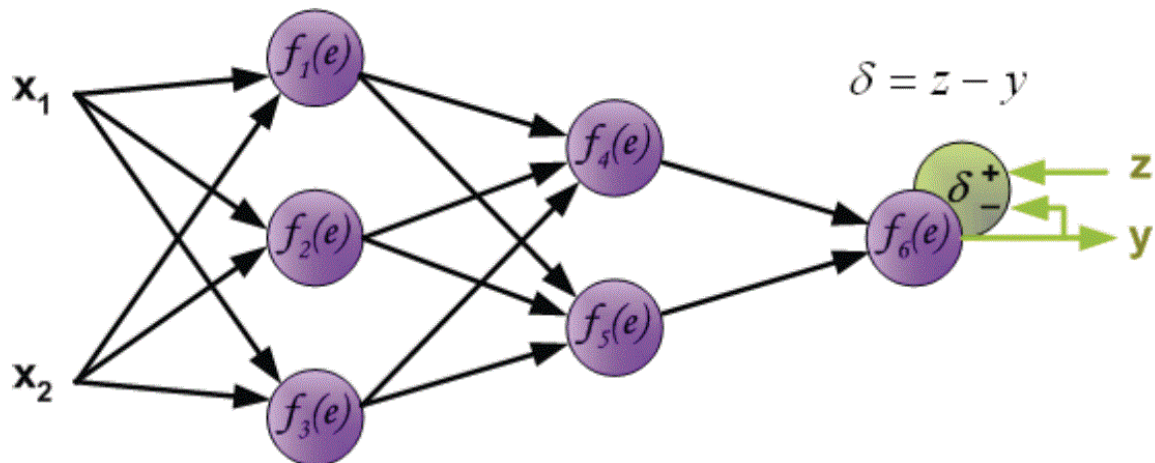
← 관심 class y
← 모든 class y

13

Basic Concepts of NN

❖ Learning: Backpropagation

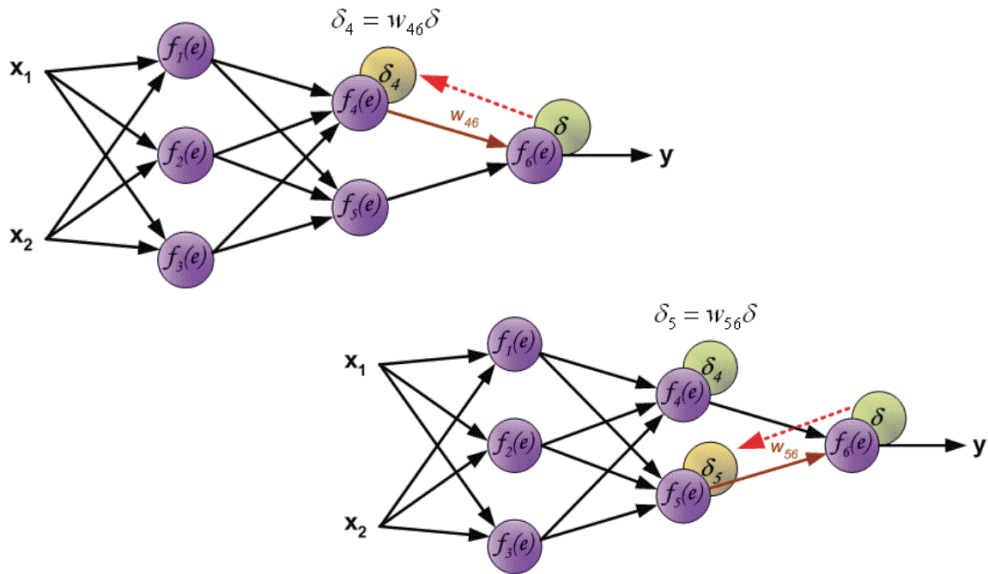
- Calculate error at the output
- Back-propagation = gradient descent + chain rule



14

Basic Concepts of NN

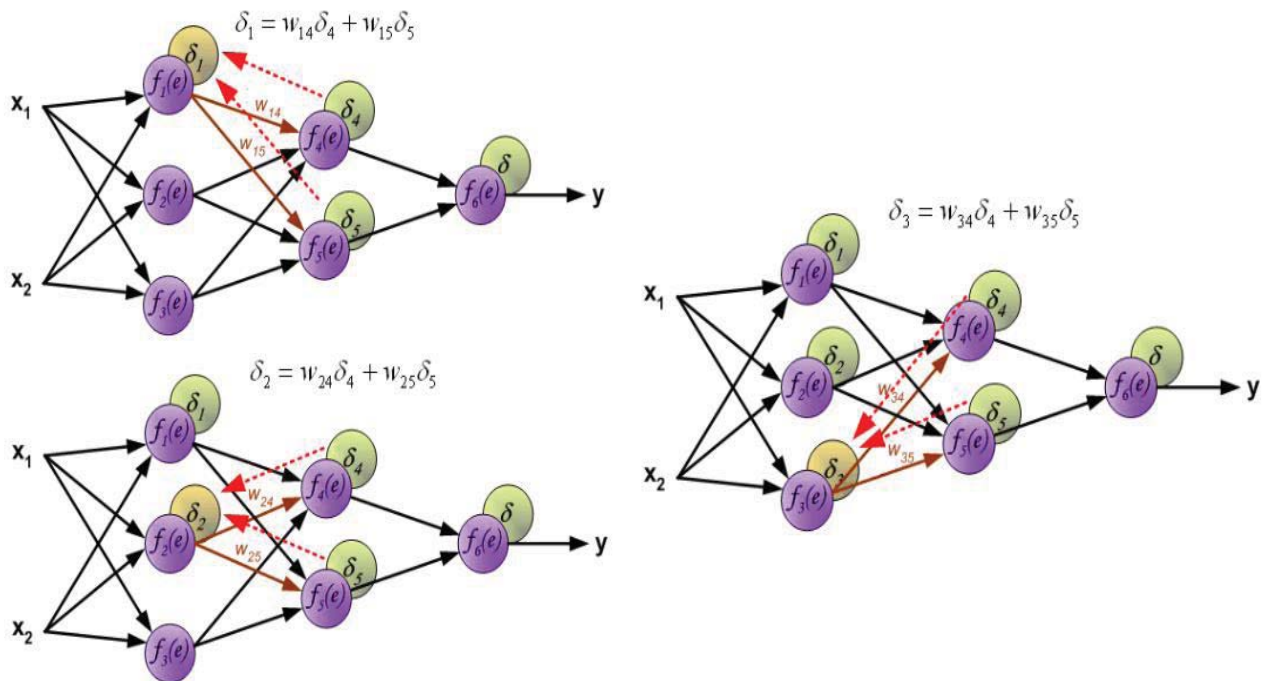
❖ Learning: Backpropagation



15

Basic Concepts of NN

❖ Learning: Backpropagation

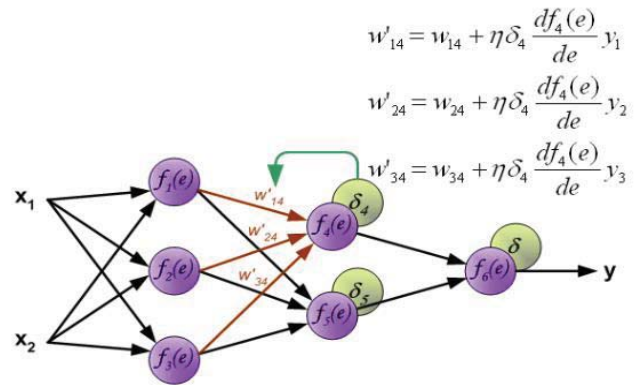
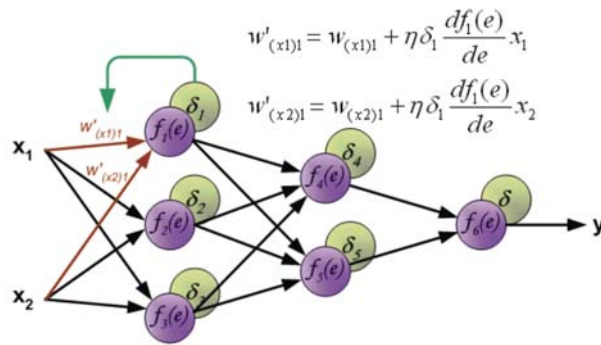


16

Basic Concepts of NN

❖ Learning: Backpropagation

- Calculate error at the output

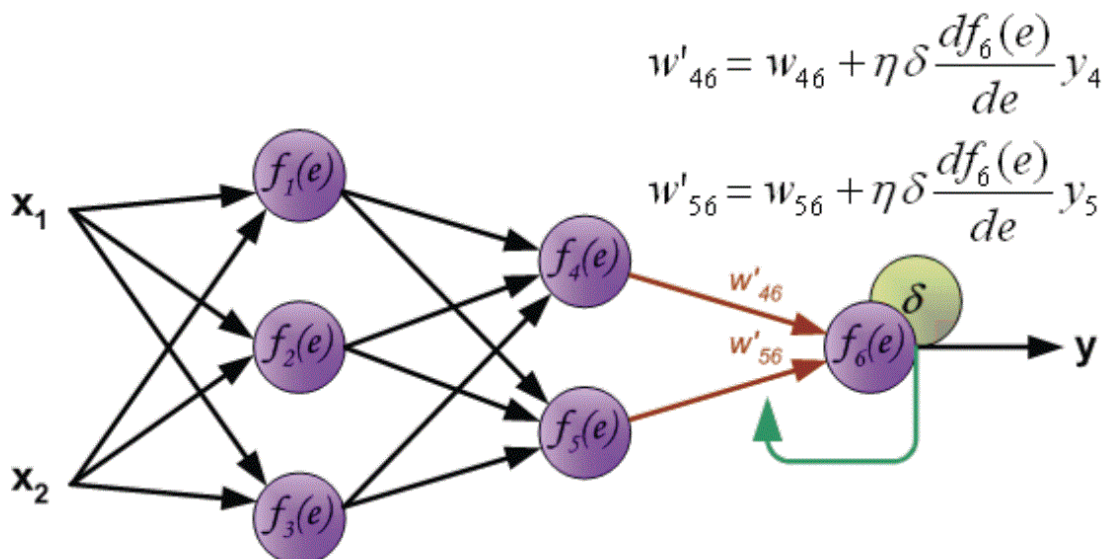


17

Basic Concepts of NN

❖ Learning: Backpropagation

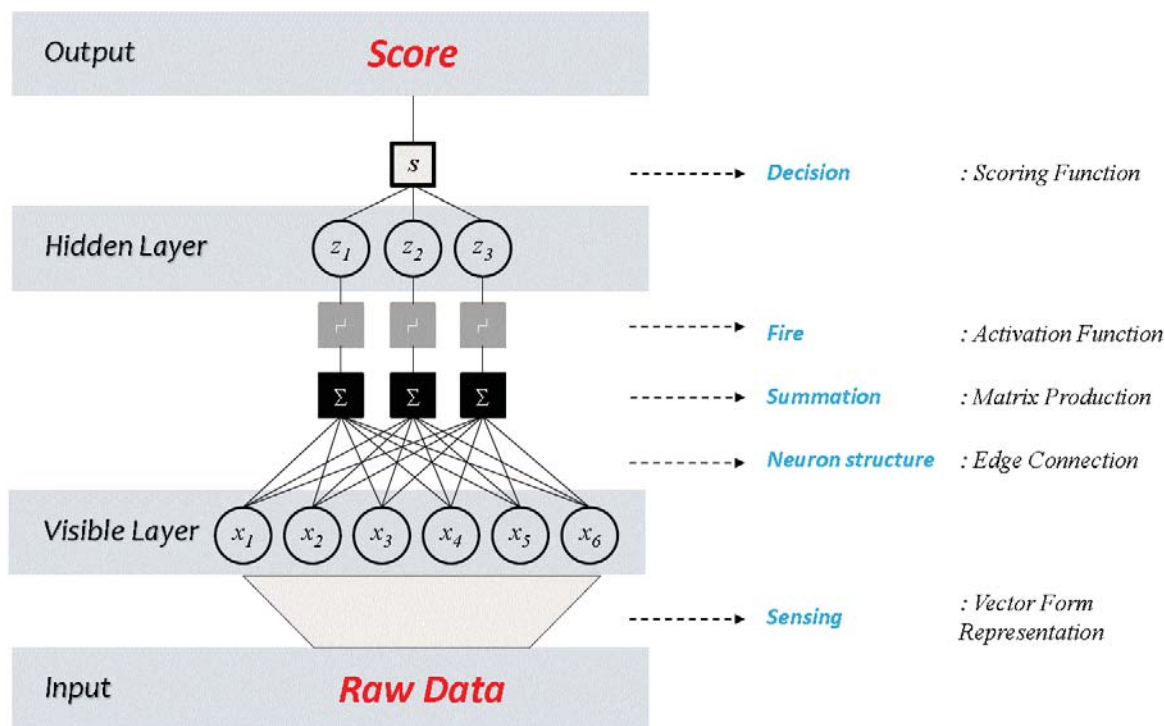
- Calculate error at the output



18

Basic Concepts of NN

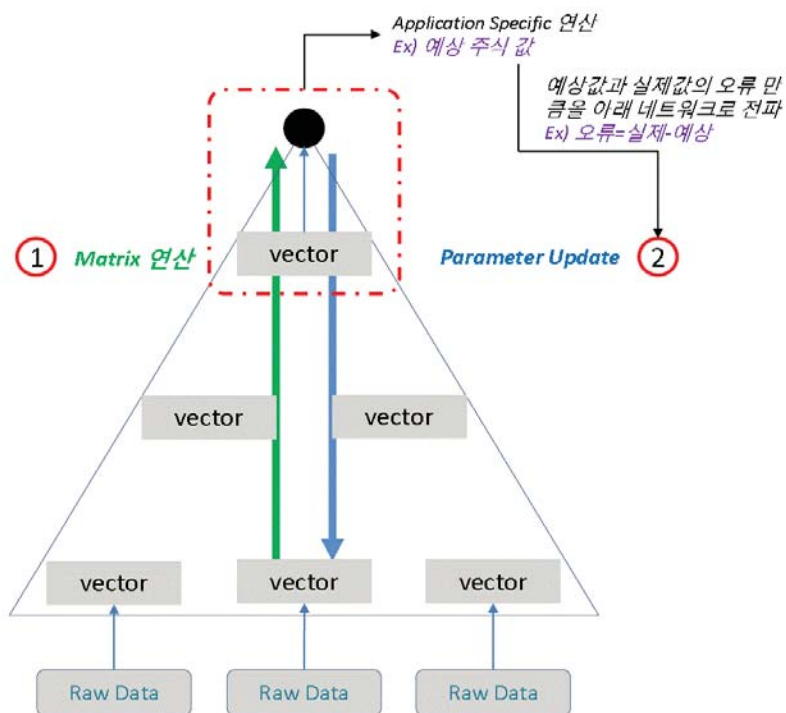
❖ Neural Network-Core Components



19

Basic Concepts of NN

❖ Neural Network-Process



20

Why? Deep Learning

❖ Why was not old NN successful?

Initialization

Local Minima

Computation Power

Data

Pre-Training

Distributed Representation

Initialization Techniques

Activation Function

Understanding ANN

Big Data

...

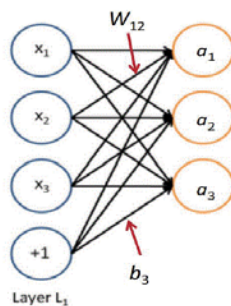
Deep Learning

21

Why? Deep Learning

❖ Neural Network-Process

- ▶ 네트워크가 깊어지고, 복잡해질수록 parameter 수가 많아짐
- ▶ Parameter가 많아질수록 Local Minima에 빠질 가능성이 높아짐



$$W = \begin{pmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} & w_{33} \end{pmatrix} \quad \mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix}$$

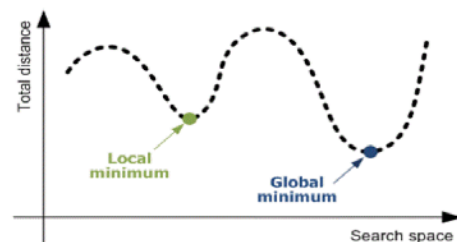
$$a_1 = f(W_{11}x_1 + W_{12}x_2 + W_{13}x_3 + b_1)$$

$$a_2 = f(W_{21}x_1 + W_{22}x_2 + W_{23}x_3 + b_2)$$

$$a_3 = f(W_{31}x_1 + W_{32}x_2 + W_{33}x_3 + b_3)$$

In Matrix notation

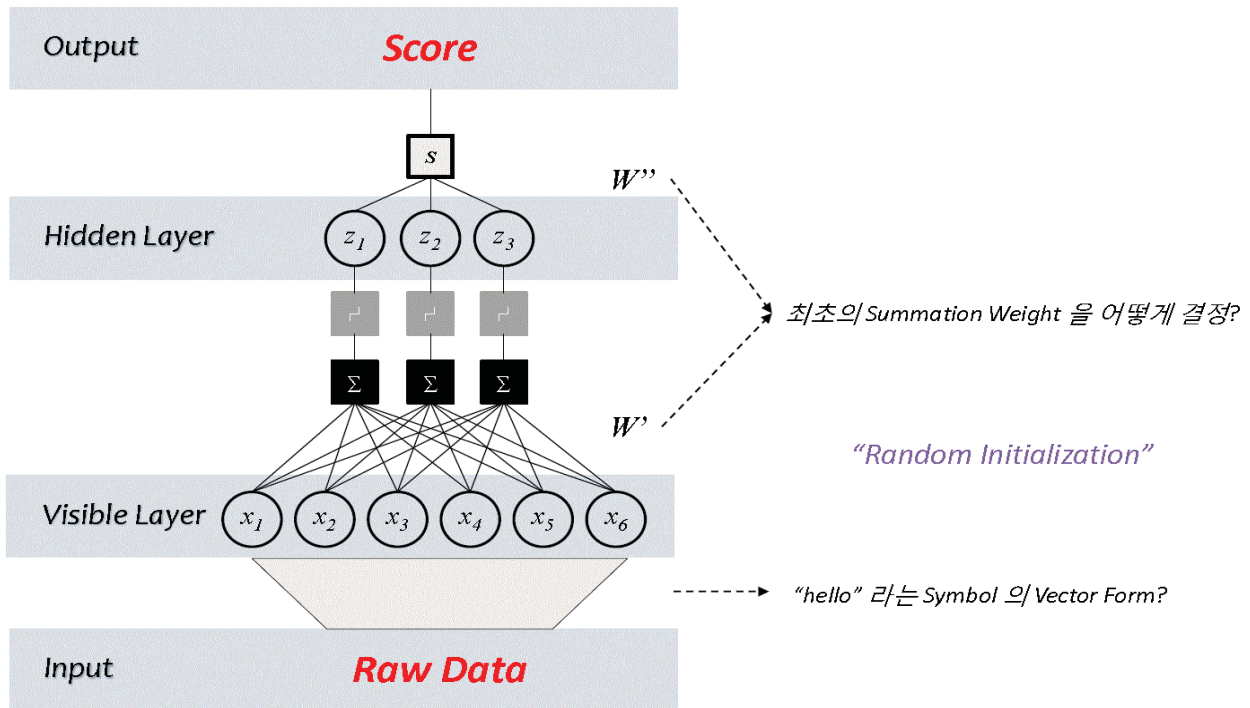
$$\mathbf{z} = \mathbf{W}\mathbf{x} + \mathbf{b}$$



22

Why? Deep Learning

❖ Initialization Problem



23

Why? Deep Learning

❖ Initialization Tip

Initial Value

초기값 Settings

- Random 하게 주되 특정 구역안에서 Random 하게 주는것이 좋다.

\tanh 를 Activation 으로 사용하는 경우

$$\text{Interval} = \left[-\sqrt{\frac{6}{fan_{in} + fan_{out}}}, \sqrt{\frac{6}{fan_{in} + fan_{out}}} \right]$$

fan_{in} = the number of units in the $(i-1)$ th layer.

fan_{out} = the number of units in the i th layer

sigmoid 를 Activation 으로 사용하는 경우

$$\text{Interval} = \left[-4\sqrt{\frac{6}{fan_{in} + fan_{out}}}, 4\sqrt{\frac{6}{fan_{in} + fan_{out}}} \right]$$

24

Why? Deep Learning

❖ Deeper Network, Harder Learning

- Network가 깊으면 깊을 수록 최종성능이 좋다. 단, 깊어질수록 Error Propagation이 어려워짐. ReLU가 사용되는 이유

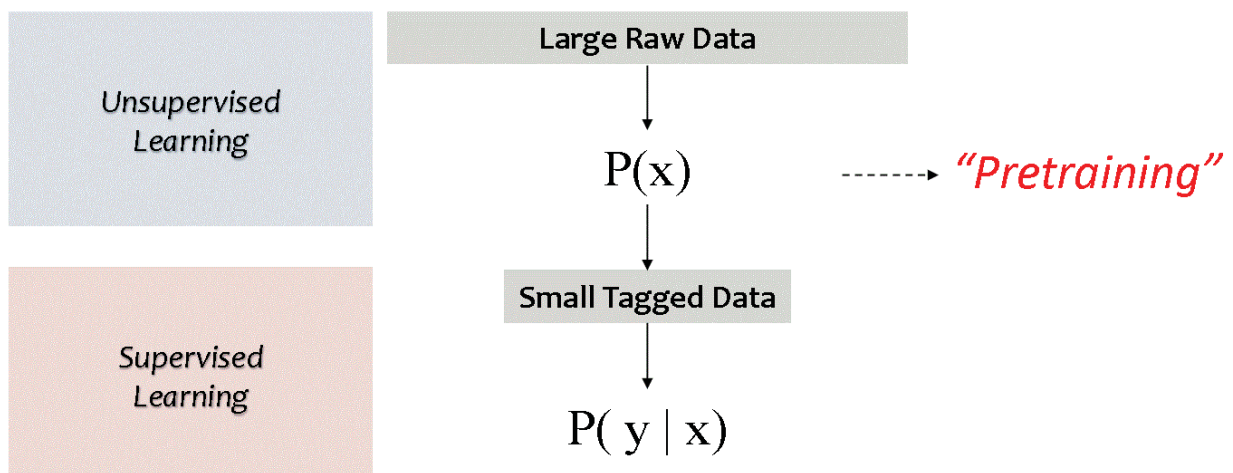


25

Why? Deep Learning

❖ Pre-Training

- Pre-training으로 NN의 성능이 비약적으로 향상됨
- AutoEncoder 계열과 Restricted Boltzmann Machine 계열이 있음

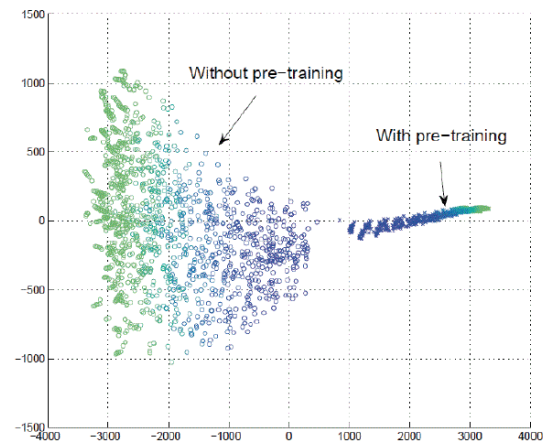


26

Why? Deep Learning

❖ Pre-Training-Performance

- Regularization hypothesis:
 - Representations good for $P(x)$ are good for $P(y|x)$
- Optimization hypothesis:
 - Unsupervised initializations start near better local minimum of supervised training error
 - Minima otherwise not achievable by random initialization



Erhan, Courville, Manzagol, Vincent, Bengio (JMLR, 2010)

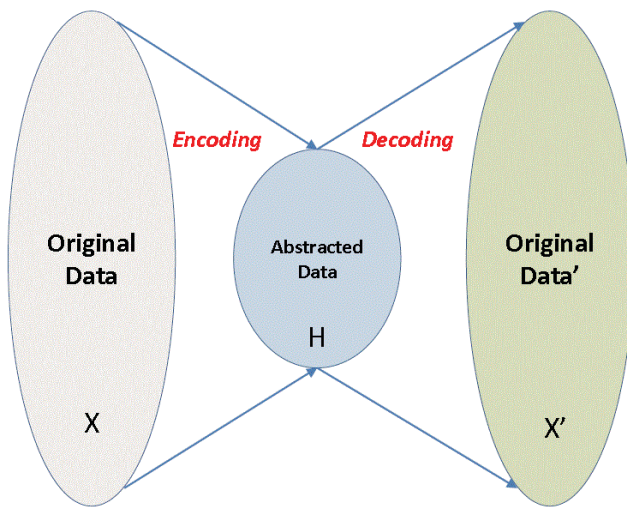
97



27

Why? Deep Learning

❖ Auto Encoder



원래의 데이터 x 를 H 로 프로젝션 시킨 후, H 로 부터 x' 를 다시 생성시킴

- 비교

- 압축 알고리즘 (Zip, MPEG, PNG ...)
- Principle Component Analysis (PCA)
- Kernel Function in SVM (original space \rightarrow hyper space)

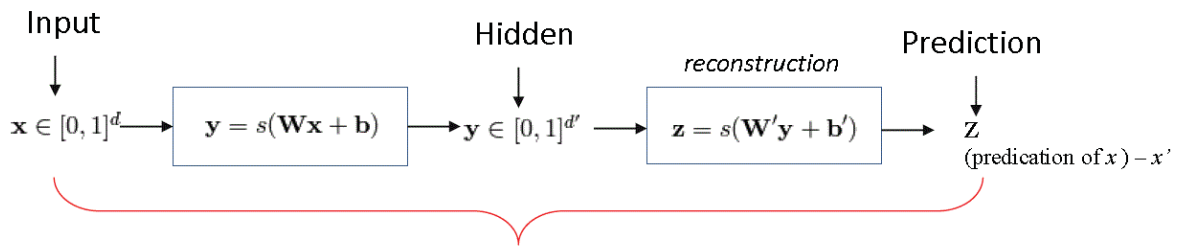
$X \rightarrow H \rightarrow X'$ 에서 $\text{Difference}(X, X')$ 가 적으면 적을수록 추상화는 완벽하게 이루어진 것이라 생각할 수 있다. 그러한 Projection 이 완벽하게 훈련된다면,

- Abstracted Data 는 그 자체로 원래 데이터를 설명하는 Feature라고 볼 수 있을 것이다.
- Feature Learning 이 자동으로 이루어지는 것이라 할 수 있음

28

Why? Deep Learning

❖ Auto Encoder



Encoding/Decoding Error

$$L(x, z) = \|x - z\|^2 \quad \dots \text{real value case}$$

$$L_H(x, z) = - \sum_{k=1}^d [x_k \log z_k + (1 - x_k) \log(1 - z_k)] \quad \dots \text{bit vector or} \\ \dots \text{vectors of bit probability}$$

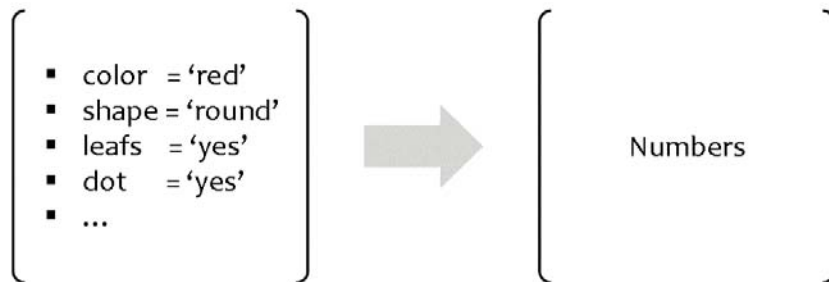
Cross-entropy

29

Learning Representation for NLP



No more handcraft feature engineering!



- 사과를 '사과'로 구별 짓는 표현방식을 스스로 학습

30

Learning Representation for NLP

❖ One-hot representation (or symbolic)

- Ex) [0 0 0 0 0 0 1 0 0 0 0 0 0 0 0 0]
- Dimensionality
 - 20K (speech) – 50K (PTB) – 500K (big vocab) – 3M (Google 1T)

❖ Continuous representation

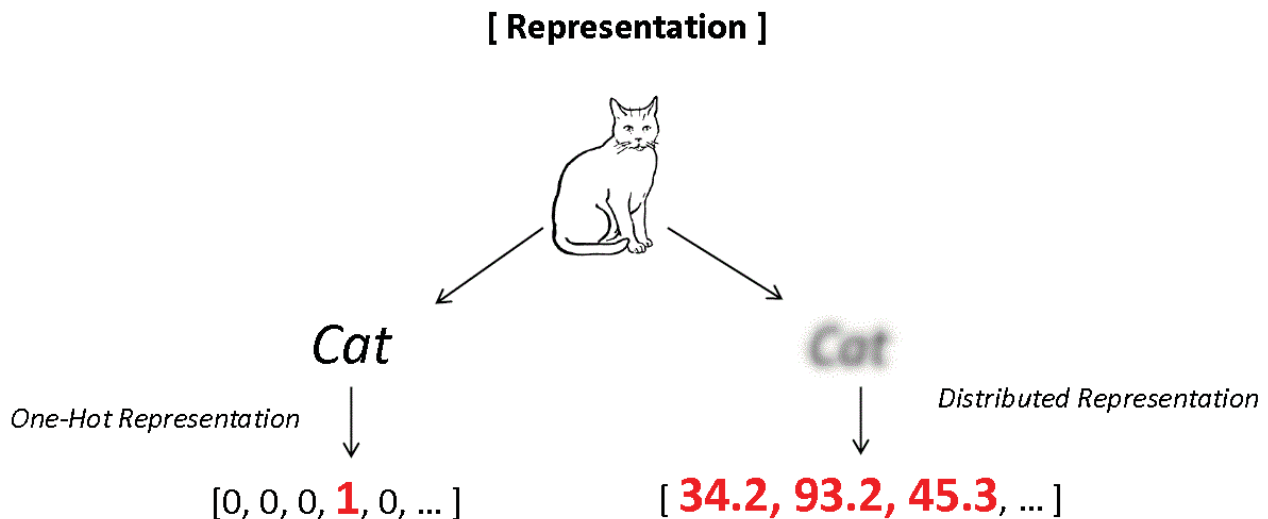
- Latent Semantic Analysis, Random projection
- Latent Dirichlet Allocation, HMM clustering
- **Distributed Representation (Neural word embedding)**
 - **Dense vector**
 - By **adding supervision** from other tasks -> **improve the representation**

31

Learning Representation for NLP

❖ Distributed Representation

- DNN이 기존 AI 방법론들에 비해 큰 의미가 있는 것은 실 세계에 있는 실제 Object를 표현할 때 Symbol에 의존하지 않는다는 점이다.

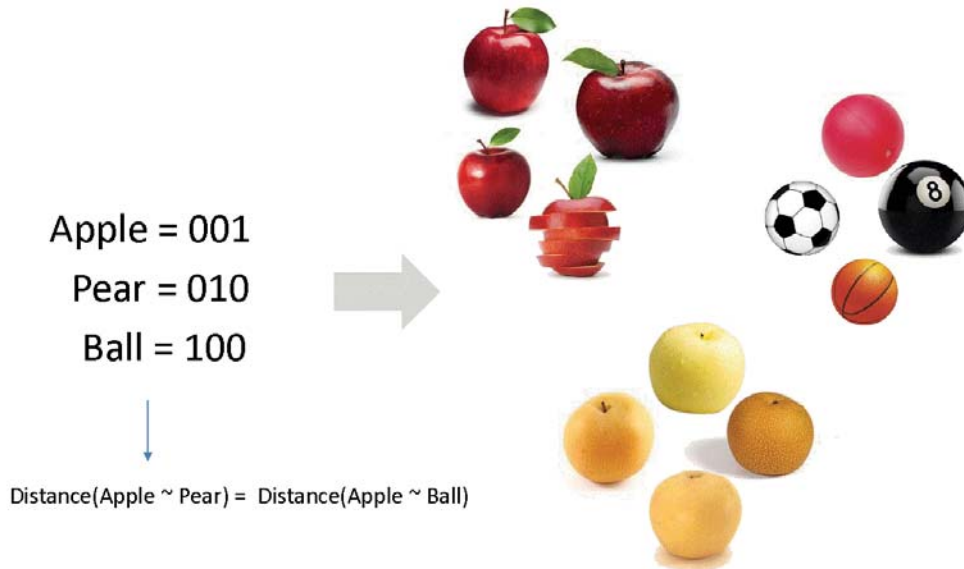


32

Learning Representation for NLP

❖ Distributed Representation

- 유사한 것은 '유사하게' 표현되어야 함
- Curse of Dimensionality 극복 가능



33

Learning Representation for NLP

Local Representation

Only one neuron (or very few) is active

Cat

[0,0,0,0,0,0,1,0,0,0,0]

- One-Hot Representation
 - **Integer Space**
 - Very Sparse
 - Very high dimensionality

Distributed Representation

many features, each of which can separately each be active or inactive

Cat

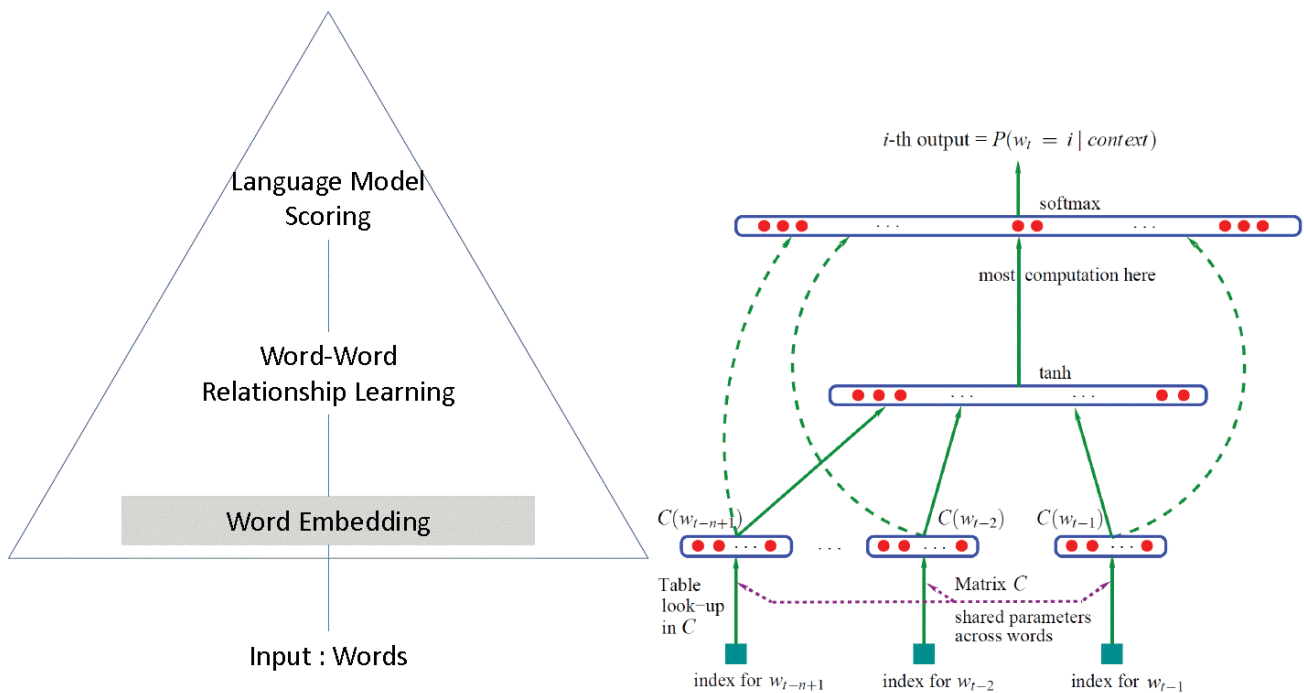
$\begin{pmatrix} -2.3 \\ 1.0 \\ 4.2 \\ 5.3 \\ 2.3 \end{pmatrix}$

- Word embedding
 - **Real value space**
 - Dense
 - Low Dimensionality

Ex) word → hash to DB Access?
It means 'integer' space.

34

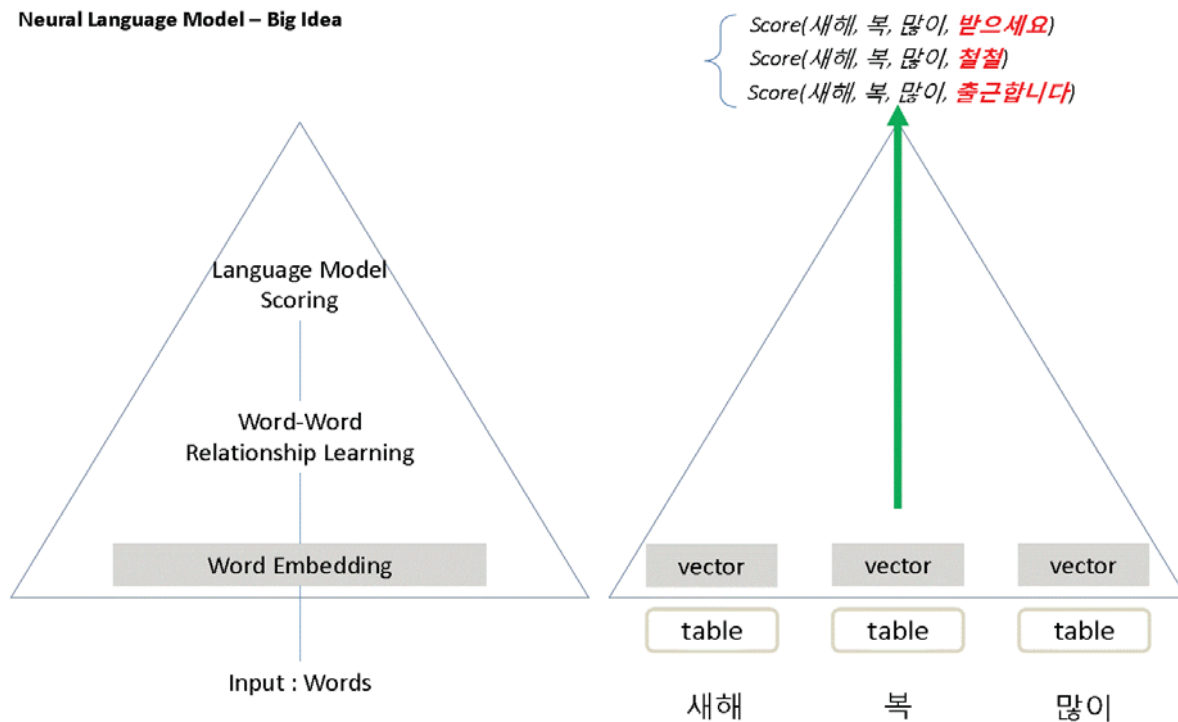
Learning Representation for NLP



35

Learning Representation for NLP

Neural Language Model – Big Idea



$Score(\text{새해}, \text{복}, \text{많이}, \text{받으세요}) > Score(\text{새해}, \text{복}, \text{많이}, \text{철철})$

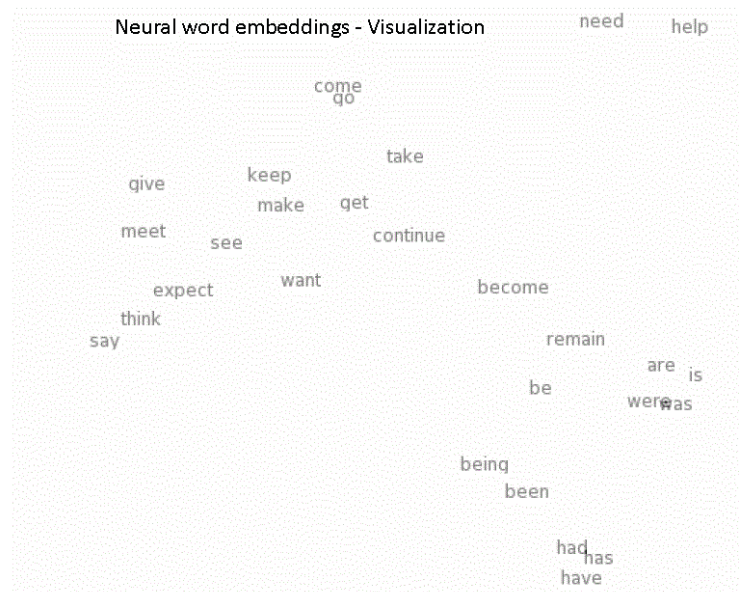
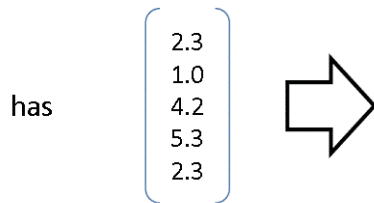
:: 실 언어데이터에 있는 패턴이 그렇지 않은 패턴보다 많이 나타나도록 NN 을 훈련시킨다.

36

Learning Representation for NLP

❖ Good One – Word Representation

We can compare words without any extra knowledge such as word net!!!

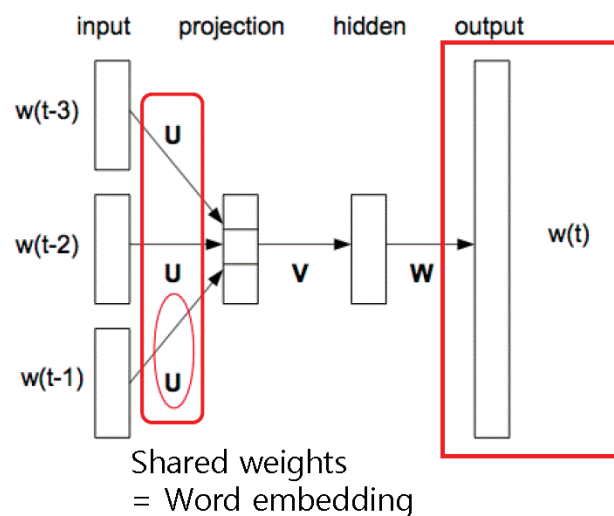


37

Learning Representation for NLP

❖ Neural Network Language Model

- Idea
 - A word and its context is a **positive** training sample
 - A random word in that same context \rightarrow **negative** training sample
 - $\text{Score}(\text{positive}) > \text{Score}(\text{neg.})$
- Training **complexity is high**
 - **Hidden layer \rightarrow output**
 - **Softmax in the output layer**
 - Hierarchical softmax
 - Negative sampling
 - Ranking(hinge loss)



Input	Dim: 1	Dim: 2	Dim: 3	Dim: 4	Dim: 5
1 (boy)	0.01	0.2	-0.04	0.05	-0.3
2 (girl)	0.02	0.22	-0.05	0.04	-0.4

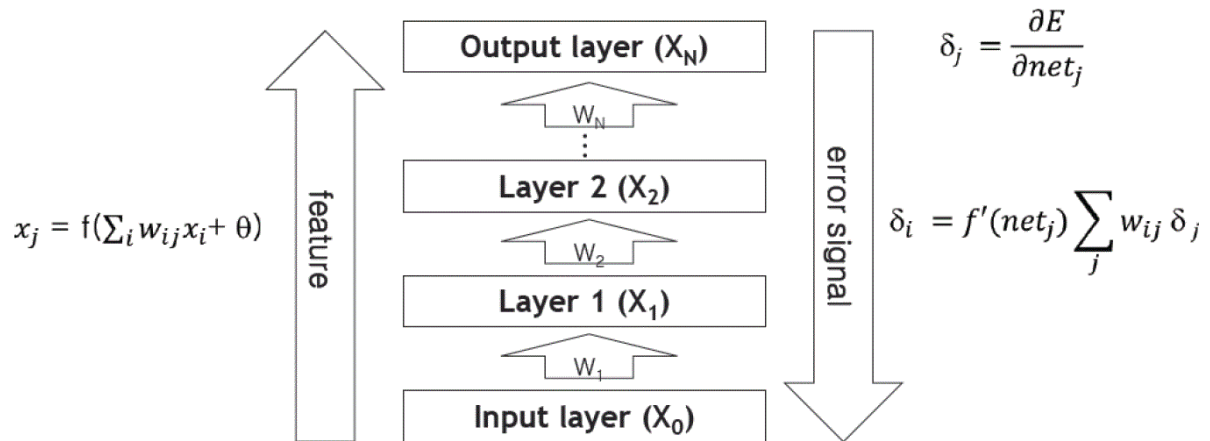
38

Learning Representation for NLP

❖ Back-Propagation Algorithm

- Gradient descent algorithm w.r.t. error E .

$$w_{ij} \leftarrow w_{ij} - \eta \frac{\partial E}{\partial w_{ij}} \qquad \frac{\partial E}{\partial w_{ij}} = \frac{\partial E}{\partial net_i} \frac{\partial net_i}{\partial w_{ij}}$$



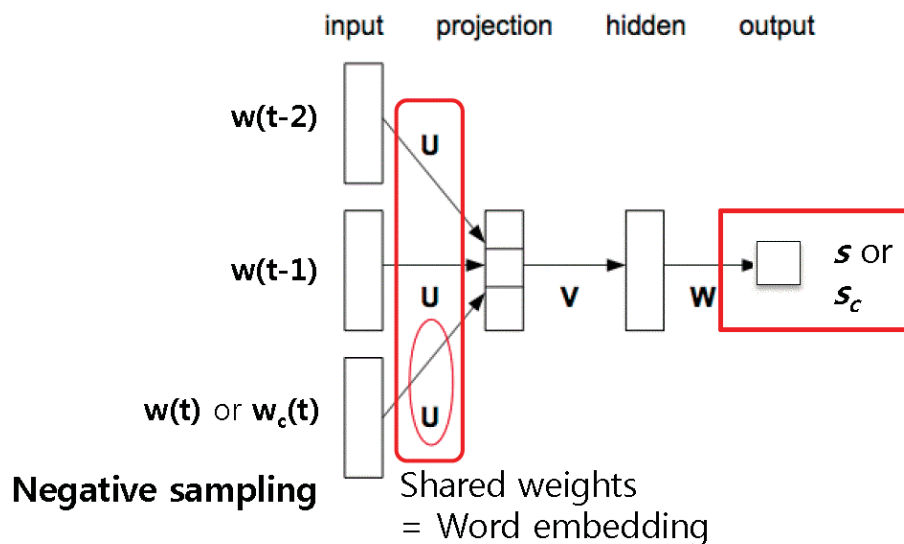
39

Learning Representation for NLP

❖ Ranking-based (Collobert)

Ranking(hinge loss)[2,11]: $\max(0, 1 - s + s_c)$

Ranking(logit loss): $\log(1 + \exp(s_c - s))$ → 본 연구 추가



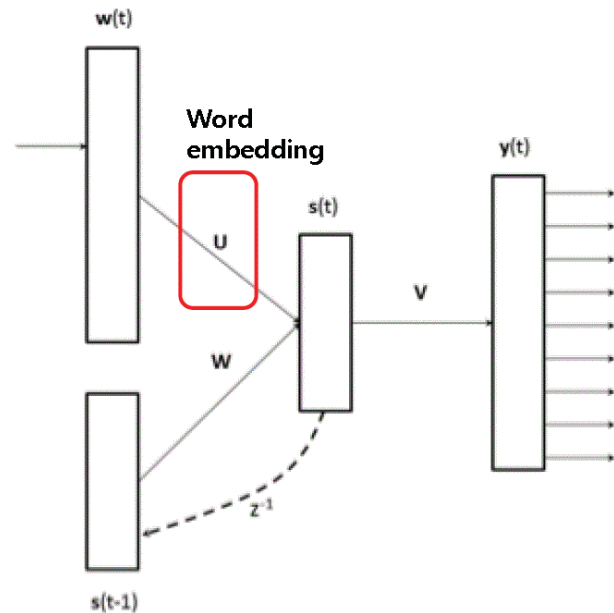
40

Learning Representation for NLP

❖ Recurrent Neural Network

• RNN

- The **hidden layer** $s(t)$ maintains a representation of the **sentence history**

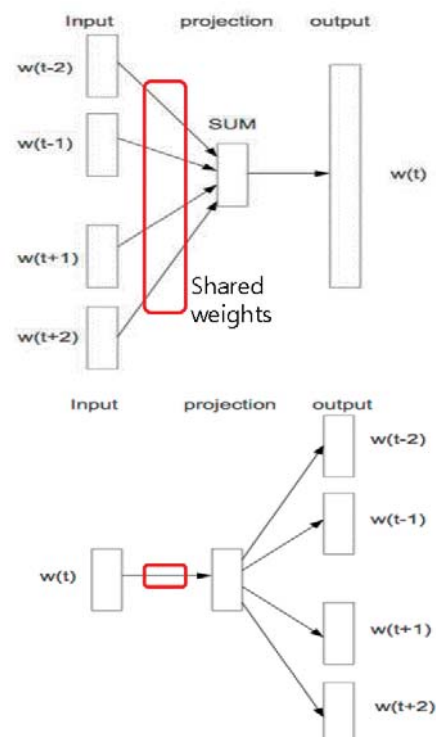


41

Learning Representation for NLP

❖ Word2Vec: CBOW, Skip-Gram

- **Remove the hidden layer → Speedup 1000x**
 - Negative sampling
 - Frequent word sampling
 - Multi-thread (no lock)
- **Continuous Bag-of-words (CBOW)**
 - Predicts the current word given the context
- **Skip-gram**
 - Predicts the surrounding words given the current word
 - **CBOW + DropOut/DropConnect**



42

Tools for Word Embedding

❖ Word2Vec

- <https://code.google.com/p/word2vec/>
- <http://deeplearning4j.org/word2vec.html#just>
- Ubutu 버전 (JAVA)
 - Googlecode에서 파일 다운로드
 - svn checkout <http://word2vec.googlecode.com/svn/trunk>
 - trunk 폴더에 파일 다운로드 됨

```
islab@islab:~$ svn checkout http://word2vec.googlecode.com/svn/trunk
A trunk/word2phrase.c
A trunk/LICENSE
A trunk/word-analogy.c
A trunk/compute-accuracy.c
A trunk/demo-analogy.sh
A trunk/demo-classes.sh
A trunk/demo-train-big-model-v1.sh
A trunk/demo-word-accuracy.sh
A trunk/demo-phrases.sh
A trunk/questions-words.txt
A trunk/demo-phrase-accuracy.sh
A trunk/demo-word.sh
A trunk/distance.c
A trunk/README.txt
A trunk/questions-phrases.txt
A trunk/word2vec.c
A trunk/makefile
체크아웃된 리비전 42.
islab@islab:~$ ls
CRF_테형
Downloads
NEMain.java
Word_Embeddings_Result
eclipse
examples.desktop
java_error_in_PYCHARM_2689.log
python
trunk
workspace
workspace_kth
workspace_lcs
workspace_twitter
```

43

Tools for Word Embedding

❖ Word2Vec 설치

- Make 명령어 실행
 - Warning은 무시

```
islab@islab:~/trunk
eclipse workspace 사진
examples.desktop workspace_kth 을악
java_error_in_PYCHARM_2689.log workspace_lcs 템플릿
python workspace_twitter
islab@islab:~$ cd trunk/
islab@islab:~/trunk$ make
gcc word2vec.c -o word2vec -lm -pthread -O3 -march=native -Wall -funroll-loops -Wno-unused-result
gcc word2phrase.c -o word2phrase -lm -pthread -O3 -march=native -Wall -funroll-loops -Wno-unused-result
gcc distance.c -o distance -lm -pthread -O3 -march=native -Wall -funroll-loops -Wno-unused-result
distance.c: In function 'main':
distance.c:31:8: warning: unused variable 'ch' [-Wunused-variable]
    char ch;
    ^
gcc word-analogy.c -o word-analogy -lm -pthread -O3 -march=native -Wall -funroll-loops -Wno-unused-result
word-analogy.c: In function 'main':
word-analogy.c:31:8: warning: unused variable 'ch' [-Wunused-variable]
    char ch;
    ^
gcc compute-accuracy.c -o compute-accuracy -lm -pthread -O3 -march=native -Wall -funroll-loops -Wno-unused-result
compute-accuracy.c: In function 'main':
compute-accuracy.c:29:109: warning: unused variable 'ch' [-Wunused-variable]
    char st1[max_size], st2[max_size], st3[max_size], st4[max_size], bestw[N][max_size], file_name[max_size], ch;
    ^
chmod +x *.sh
islab@islab:~/trunk$
```

44

Tools for Word Embedding

❖ Word2Vec 실행

```
islab@islab: ~/word2vec/source
islab@islab:~/word2vec/source$ ls
LICENSE                demo-train-big-model-v1.sh  questions-words.txt
README.txt             demo-word-accuracy.sh      text8
compute-accuracy      demo-word.sh               word-analogy
compute-accuracy.c    distance                   word-analogy.c
demo-analogy.sh       distance.c                 word2phrase
demo-classes.sh       makefile                  word2phrase.c
demo-phrase-accuracy.sh naver_word2vec.txt        word2vec
demo-phrases.sh       questions-phrases.txt     word2vec.c
islab@islab:~/word2vec/source$ ./word2vec -train ../data/yahoo_wordembedding.txt -o
output ../data/naver_word2vec.txt -size 64 -iter 50
Starting training using file ../data/yahoo_wordembedding.txt
Vocab size: 338381
Words in train file: 247197230
Alpha: 0.049916 Progress: 0.17% Words/thread/sec: 474.26k
```

45

Tools for Word Embedding

❖ Word2Vec parameters

- -output
 - 출력파일
- -size
 - 생성할 word vector의 차원 (default value: 100)
- -windows
 - Max skip length (default value: 5)
- -cbow
 - 1: continuous bag of word model, 0: skip-gram model
- -iter
 - 반복 횟수 (default value: 5)
- -min-count
 - 지정 횟수보다 작게 나온 단어 제거 (default value: 5)
- -save-vocab
 - 단어 목록 출력 파일
- -read-vocab
 - 단어목록을 미리 지정

46

Tools for Word Embedding

❖ Word2Vec 학습파일 포맷

- -train
- 한 문장 별로 한 라인에 문장 자질로 구성

지금/MAG 도/JX 경기도/NNP 나/NP 강원도/NNP 산세/NNG 종/VA 은/ETM 곳/NNG 예/JKB 가면/NNG 산/NNG 과/JC 숲/NNG 을/JKO 온통/MAG 파헤치/VV 는/ETM 골프장/NN
달팽/XR 하/XSA ㄴ/ETM 산/NNG 을/JKO 현/VV 어/EC 내/VX 고/EC 언덕/NNG 을/JKO 만드/NA 러니/EC 얼마나/MAG 힘/NNG 이/JKS 들/VV 쟁/EP 는/ETM 가/NNG ./SF
케다가/MAG 그/MM 위/NNB 예/JKB 잔디/NNG 까지/JX 입혀/VV 러니/EC 문제/NNG 가/JKS 한/NNP 들이/NNG 아니/VCM 다/EF ./SF
원래/NNG 우리나라/NNG 기후/NNG 예/JKB 는/JX 잔디/NNG 가/JKS 잘/MAG 맞/VV 지/EC 않/VX 는다/EF ./SF
야산/NNG 예/JKB 잡목/NNG 과/JC 덩굴/NNG 숲/NNG 만/JX 우거졌/NA 지/VX 초원/NNG 이/JKS 생겨나/VV 지/EC 않/VX 은/ETM 첫/NNB 도/JX 기/NNG 후/NNG 때문/NNB

➤ Tutorial

- <http://alexminnaar.com/word2vec-tutorial-part-ii-the-continuous-bag-of-words-model.html>

Tools for Word Embedding

❖ Ranking-based Model

- <https://bitbucket.org/aboSamoor/word2embeddings>
- Python, Theano
- Result file (pickle)

- `a = np.load(sys.argv[2])`

References

- ❖ Ronan Collbert, et al. “Natural Language Processing (Almost) from Scratch,” *Journal of Machine Learning Research*, 2011.
- ❖ Mikolov, T., et al. “Recurrent Neural Network based Language Model,” 2010.
- ❖ Mikolov, T., et al., “Distributed Representations of Words and Phrases and their Compositionality,” *NIPS*, 2013.
- ❖ Le, Q., Mikolov, T., “Distributed Representations of Sentences and Documents,” *ICML*, 2014.
- ❖ Kalchbrenner, N., Grefenstette, E. and Blunsom, P. “A Convolutional Neural Network for Modelling Sentences,” *ACL*, 2014.
- ❖ Kim, Y. “Convolutional Neural Networks for Sentence Classification,” *EMNLP*, 2014.
- ❖ Al-Rfou, R., Perozzi, B., Skiena, S., “Polyglot: Distributed Word Representations for Multilingual NLP,” *ACL*, 2013.
- ❖ 정상근, “Introduction to Deep Learning,” *자연언어처리 튜토리얼*, 2015.
- ❖ 이창기, “Word and Phrase Embedding,” *자연언어처리 튜토리얼*, 2015.

49

Thank you for your attention!

<http://web.donga.ac.kr/yjko/>

고영중