

Text Processing

(good survey at https://github.com/kk7nc/Text_Classification)

2023. 8

Yongjin Jeong, KwangWoon University

[참고] 본 자료에는 인터넷에서 다운받아 사용한 그림이나 수식들이 일부
있으니 다른 용도로 사용하거나 외부로 유출을 금해 주시기 바랍니다.

Text Processing

- **Applications**

- Language translation
- Sentiment analysis (positive, negative, neutral) (기쁨, 우울, 슬픔, 등)
- Spam filtering (emails, messages)
- Information extraction (from SNS, Web, Documents)
- AI Speaker, Chatbot, Question answering

- **Language**

- **Token**: basic text processing unit (character, **word**, n-gram, POS(part-of-speech, 품사), syllable(음절:자음+모음))
- **Corpus** (말뭉치): collection of text (or **document**) dataset
- **Stop words** (불용어): most common words in a language ("the", "is", etc)

- **Tools**

- NLTK (Natural Language ToolKit): all sorts of tasks from tokenization, stemming, tagging, parsing, and beyond
- KoNLPy (NLP library for Korean Language)
- Gensim.Word2Vec(), tf.Keras.Embedding(), CountVectorizer(), BeautifulSoup()

Text Processing

- **Tokenization**

- way of separating a piece of text into smaller units called tokens.
- Tokens can be either words, characters, partial words (n-gram characters)

- **NLTK (Natural Language Toolkit)**

```
1 from nltk.tokenize import word_tokenize
2
3 text = "After sleeping for four hours, he decided to sleep for another four"
4 tokens = word_tokenize(text)
5 print(tokens)
```

```
['After', 'sleeping', 'for', 'four', 'hours', ',', 'he', 'decided', 'to', 'sleep', 'for', 'another', 'four']
```

```
1 example_sent = "This is a sample sentence, showing off the stop words filtration."
2 stop_words = set(stopwords.words('english'))
3 word_tokens = word_tokenize(example_sent)
4 filtered_sentence = [w for w in word_tokens if not w in stop_words]
5 print(word_tokens)
6 print(filtered_sentence)
```

```
['This', 'is', 'a', 'sample', 'sentence', ',', 'showing', 'off', 'the', 'stop', 'words', 'filtration', '.']
['This', 'sample', 'sentence', ',', 'showing', 'stop', 'words', 'filtration', '.']
```

Text Processing

- KoNLPy (Korean Natural Language Processing for Python: 코엔엘파이)
 - Morpheme tokenizer (형태소 분석기):
 - Okt(Open Korea Text, 옛이름: Twitter)
 - Mecab(메캅), Komoran(코모란), Hannanum(한나눔)
 - Kkma (꼬꼬마)

```
1 from konlpy.tag import Okt
2 okt=Okt()
3
4 print(okt.morphs("우리가 이 과제를 잘 할 수 있을까?"))
5 print(okt.pos("우리가 이 과제를 잘 할 수 있을까?", norm=True, stem=True))
6 print(okt.nouns("우리가 이 과제를 잘 할 수 있을까?"))
```

```
[ '우리', '가', '이', '과제', '를', '잘', '할', '수', '있을까', '?' ]
[ ('우리', 'Noun'), ('가', 'Josa'), ('이', 'Noun'), ('과제', 'Noun'), ('를', 'Josa'), ('자다', 'Verb'), ('하다', 'Verb'), (
['우리', '이', '과제', '수']
```

```
1 from konlpy.tag import Okt
2 okt=Okt()
3 word_tags = okt.pos("우리가 이 과제를 잘 할 수 있을까?", norm=True, stem=True)
4 print(word_tags)
5 stop_words = [word[0] for word in word_tags if word[1]=="Josa"]
6 print (stop_words)
```

```
[('우리', 'Noun'), ('가', 'Josa'), ('이', 'Noun'), ('과제', 'Noun'), ('를', 'Josa'), ('자다', 'Verb'), ('하다', 'Verb'), (
['가', '를']
```

Text Processing - steps

- **Data preprocessing:**
 - **Tokenization:** convert sentences to words
 - **Removing punctuations and tags**
 - **Removing stop words**
 - **Stemming:** words are reduced to a root (ex: formalize -> formal, allowance -> allow)
 - **Lemmatization** (ex: is, are, am -> be, having -> have)
- **Feature Extraction**
 - One-hot encoding
 - **BOW (Bag of Words):** each sentence is represented with counts (or presence) of words (discards word order, hence **ignoring the context and meaning** in the document), **sparse**
 - `CountVectorizer()`, `TfidfVectorizer()`
 - **Word Embedding:** represent each word with a **meaningful n-dimensional vector**, **dense**
 - `Word2vec()`: put words of similar meaning (or context) to closer places
 - `tf.keras.Embedding()`: vector representation is **trained** (only used as the first layer in deep learning)
- **Choose ML algorithms**
 - Classical ML (simple text classification) or Deep Learning (sentiment analysis or language translation)

Text Processing

- **BoW** (Bag of Words)

- From http://www.datasciencecourse.org/notes/free_text/
 - Doc1 = "The goal of this lecture is to explain the basics of free text processing"
 - Doc2 = "The bag of words model is one such approach"
 - Doc3 = "Text processing via bag of words"

$$X = \begin{array}{c} \begin{array}{cccccccccc} \text{the} & \text{is} & \text{of} & \text{goal} & \text{lecture} & \text{bag} & \text{words} & \text{via} & \text{text} & \text{approach} \\ 2 & 1 & 2 & 1 & 1 & 0 & 0 & 0 & 1 & 0 \end{array} \\ \begin{array}{l} \text{Document 1} \\ \text{Document 2} \\ \text{Document 3} \end{array} \end{array}$$



$$X = \begin{array}{c} \begin{array}{cccc} \text{the} & \text{is} & \text{of} & \text{goal} \\ 0.8 & 0.4 & 0 & 1.1 \\ 0.4 & 0.4 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{array} \end{array}$$

$$\text{idf}_j = \log \left(\frac{\# \text{ documents}}{\# \text{ documents with word } j} \right)$$

$$\text{idf}_{\text{of}} = \log \left(\frac{3}{3} \right) = 0$$

$$\text{idf}_{\text{is}} = \log \left(\frac{3}{2} \right) = 0.405$$

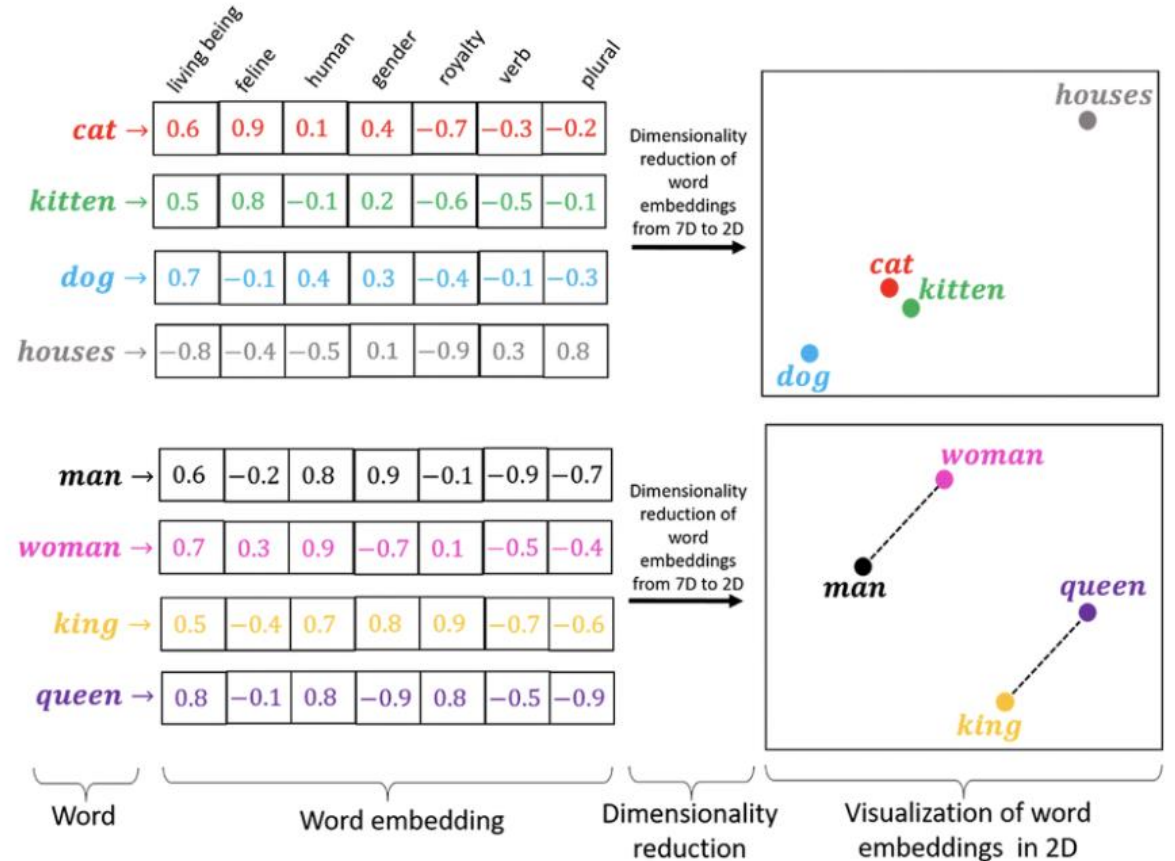
$$\text{idf}_{\text{goal}} = \log \left(\frac{3}{1} \right) = 1.098$$

Text Processing

- **Word Embedding**

- Every word has a unique word embedding (or "vector"), and similar words end up with similar embedding values.
- [Word2vec](#) from google: based on proximity
- [Glove](#): pre-trained (from Stanford)
- [Embedding\(\)](#) layer in deep learning: trained for the specific purpose

```
1 from tensorflow.keras.models import Sequential
2 from tensorflow.keras.layers import SimpleRNN, Dense, LSTM, GRU
3
4 model = Sequential()
5 model.add(Embedding(input_dim=10000, output_dim=8,
6 | | | | | input_length=maxlen))
7 model.add(LSTM(128))
8 model.add(Flatten())
9 model.add(Dense(1, activation='sigmoid'))
```



One-Hot encoding

- **One-Hot encoding**
 - Each word is written or encoded as one hot vector, with each one hot vector being **unique**.
 - One word is represented as a **vector**, therefore a sentence is represented as an array of vectors or a **matrix**.
 - A list of sentences will end up with a **three dimensional tensor** which can be fed to the Neural network.

The cat sat on the mat

The: [0 1 0 0 0 0]

cat: [0 0 1 0 0 0]

sat: [0 0 0 1 0 0]

on: [0 0 0 0 1 0]

the: [0 0 0 0 0 1]

mat: [0 0 0 0 0 0 1]

Bag of Words (BoW)

- A representation of text that describes the **occurrence of words** within a document
- Just **keep track of word counts** and **disregard the grammatical details** (word order), that's why called “bag” of words
- Occurrence of words:
 - **Counts**
 - **Frequencies**
- Limitations
 - **Sparsity**: little information with large representation space
 - **Vocabulary**: should be carefully designed due to impact on sparsity
 - **Meaning**: context and meaning of words (semantics) are disregarded
- Need text cleaning techniques
 - Ignoring cases, punctuations, frequent words (“a”, “of”, etc)
 - Fixing misspelled words
 - Reducing words to their stem (e.g. “play” from “playing”): stemization

Bag of Words (BoW)

- **N-gram**
 - N-token sequence of words
 - For example, 2-gram (more commonly called a bigram) is a two-word sequence of words (e.g. “please turn”, “turn your”, or “your homework”)
 - Often a simple bigram approach is better than a 1-gram bag-of-words model for tasks like documentation classification.
- **CountVectorizer()**
 - Document Term matrix
- **TfidfVectorizer()**
 - Term frequency (tf) : frequency of the words
 - Inverse Document frequency (Idf) : how rare the word is across documents
- **Keras**
 - `Keras.preprocessing.text.Tokenizer.fit_on_texts()`
 - `Keras.preprocessing.sequence.pad_sequences()`

Document Term Matrix

- From http://www.datasciencecourse.org/notes/free_text/
 - Doc1 = "The goal of this lecture is to explain the basics of free text processing"
 - Doc2 = "The bag of words model is one such approach"
 - Doc3 = "Text processing via bag of words"

$$X = \begin{matrix} & \begin{matrix} \text{the} & \text{is} & \text{of} & \text{goal} & \text{lecture} & \text{bag} & \text{words} & \text{via} & \text{text} & \text{approach} \end{matrix} \\ \begin{bmatrix} 2 & 1 & 2 & 1 & 1 & 0 & 0 & 0 & 1 & 0 \\ 1 & 1 & 1 & 0 & 0 & 1 & 1 & 0 & 0 & \dots & 1 \\ 0 & 0 & 1 & 0 & 0 & 1 & 1 & 1 & 1 & 0 \end{bmatrix} & \begin{matrix} \text{Document 1} \\ \text{Document 2} \\ \text{Document 3} \end{matrix} \end{matrix}$$

TFIDF

- **Term frequency**
 - Counts of each word in a document
 - $tf_{i,j}$ = frequency of word j in document i
- **Inverse document frequency**
 - Term frequencies tend to be “overloaded” with very common words (“the”, “is”, “of”, etc)
 - Idea if inverse document frequency weight words negatively in proportion to how often they occur in the entire set of documents

$$idf_j = \log \left(\frac{\# \text{ documents}}{\# \text{ documents with word } j} \right)$$

TFIDF

$$X = \begin{matrix} & \text{the} & \text{is} & \text{of} & \text{goal} & \text{lecture} & \text{bag} & \text{words} & \text{via} & \text{text} & \text{approach} \\ \begin{matrix} \text{Document 1} \\ \text{Document 2} \\ \text{Document 3} \end{matrix} & \begin{bmatrix} 2 & 1 & 2 & 1 & 1 & 0 & 0 & 0 & 1 & 0 \\ 1 & 1 & 1 & 0 & 0 & 1 & 1 & 0 & 0 & \dots & 1 \\ 0 & 0 & 1 & 0 & 0 & 1 & 1 & 1 & 1 & 1 & 0 \end{bmatrix} \end{matrix} \quad \rightarrow \quad X = \begin{bmatrix} 0.8 & 0.4 & 0 & 1.1 \\ 0.4 & 0.4 & 0 & 0 & \dots \\ 0 & 0 & 0 & 0 & \dots \end{bmatrix}$$

$$\text{idf}_{\text{of}} = \log \left(\frac{3}{3} \right) = 0$$

$$\text{idf}_{\text{is}} = \log \left(\frac{3}{2} \right) = 0.405$$

$$\text{idf}_{\text{goal}} = \log \left(\frac{3}{1} \right) = 1.098$$

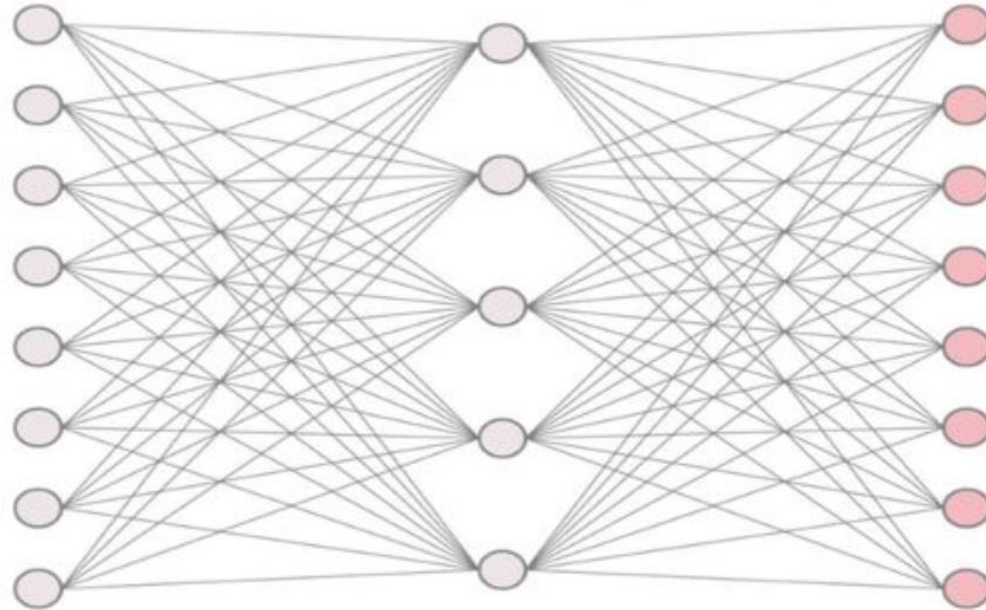
[note] In `TfidfVectorizer()`, if `smooth_idf=true` (default), the constant "1" is added to the numerator and denominator of the idf as if an extra document was seen containing every term in the collection exactly once, which prevents zero divisions: $\text{idf}(t) = \log \left[\frac{(1 + n)}{(1 + \text{df}(t))} \right] + 1$

Word Embedding

- Techniques to compute Word embedding
 - Using Supervised learning
 - Take an NLP problem and try to solve it. In that pursuit as a [side effect](#), you get word embedding.
 - Keras Embedding() layer
 - Using Self-supervised learning
 - Word2Vec (more popular)
 - Glove

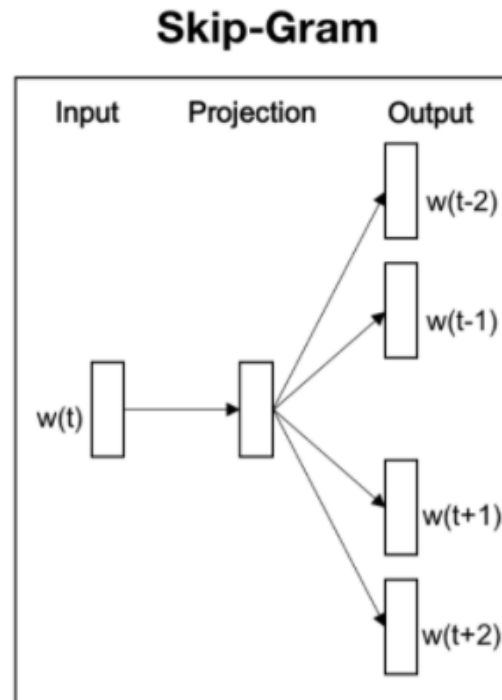
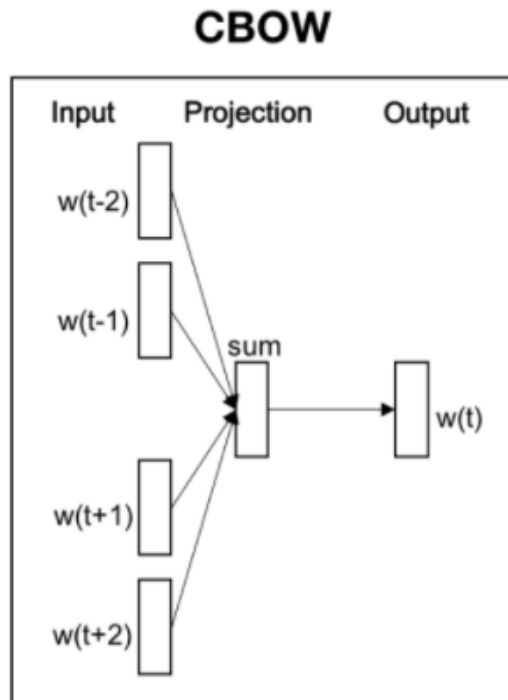
Word2Vec

- Predict words using context
 - Concept: '비슷한 위치 (근처)에서 등장하는 단어들은 비슷한 의미를 가진다'
- Word2Vec Neural Network model
 - NN with single hidden layer
 - Often used for auto-encoder to compress input vector in hidden layer



Word2Vec

- **Two versions: CBOW(continuous bag of words) and Skip-Gram**
 - CBOW: 문맥 단어를 보고 기준 단어가 무엇인지 예측하는 모델
 - Skip-Gram: 기준 단어를 보고 어떤 문맥 단어가 등장할지 예측하는 모델 (in most cases, skip-gram is better)



Word2Vec - CBOW

- (ex) "I like playing football with my friends" with window = 2

center word context words

I like playing football with my friends

I like playing football with my friends

I like playing football with my friends

I like playing football with my friends

I like playing football with my friends

I like playing football with my friends

I like playing football with my friends

y	X
center word	context words
[1,0,0,0,0,0,0]	[0,1,0,0,0,0,0] [0,0,1,0,0,0,0]
[0,1,0,0,0,0,0]	[1,0,0,0,0,0,0] [0,0,1,0,0,0,0] [0,0,0,1,0,0,0]
[0,0,1,0,0,0,0]	[1,0,0,0,0,0,0] [0,1,0,0,0,0,0] [0,0,0,1,0,0,0] [0,0,0,0,1,0,0]
[0,0,0,1,0,0,0]	[0,1,0,0,0,0,0] [0,0,1,0,0,0,0] [0,0,0,0,1,0,0] [0,0,0,0,0,1,0]
[0,0,0,0,1,0,0]	[0,0,1,0,0,0,0] [0,0,0,1,0,0,0] [0,0,0,0,0,1,0] [0,0,0,0,0,0,1]
[0,0,0,0,0,1,0]	[1,0,0,1,0,0,0] [0,0,0,0,1,0,0] [0,0,0,0,0,0,1]
[0,0,0,0,0,0,1]	[0,0,0,0,1,0,0] [0,0,0,0,0,1,0]

Training samples

(I, like)
(I, playing)

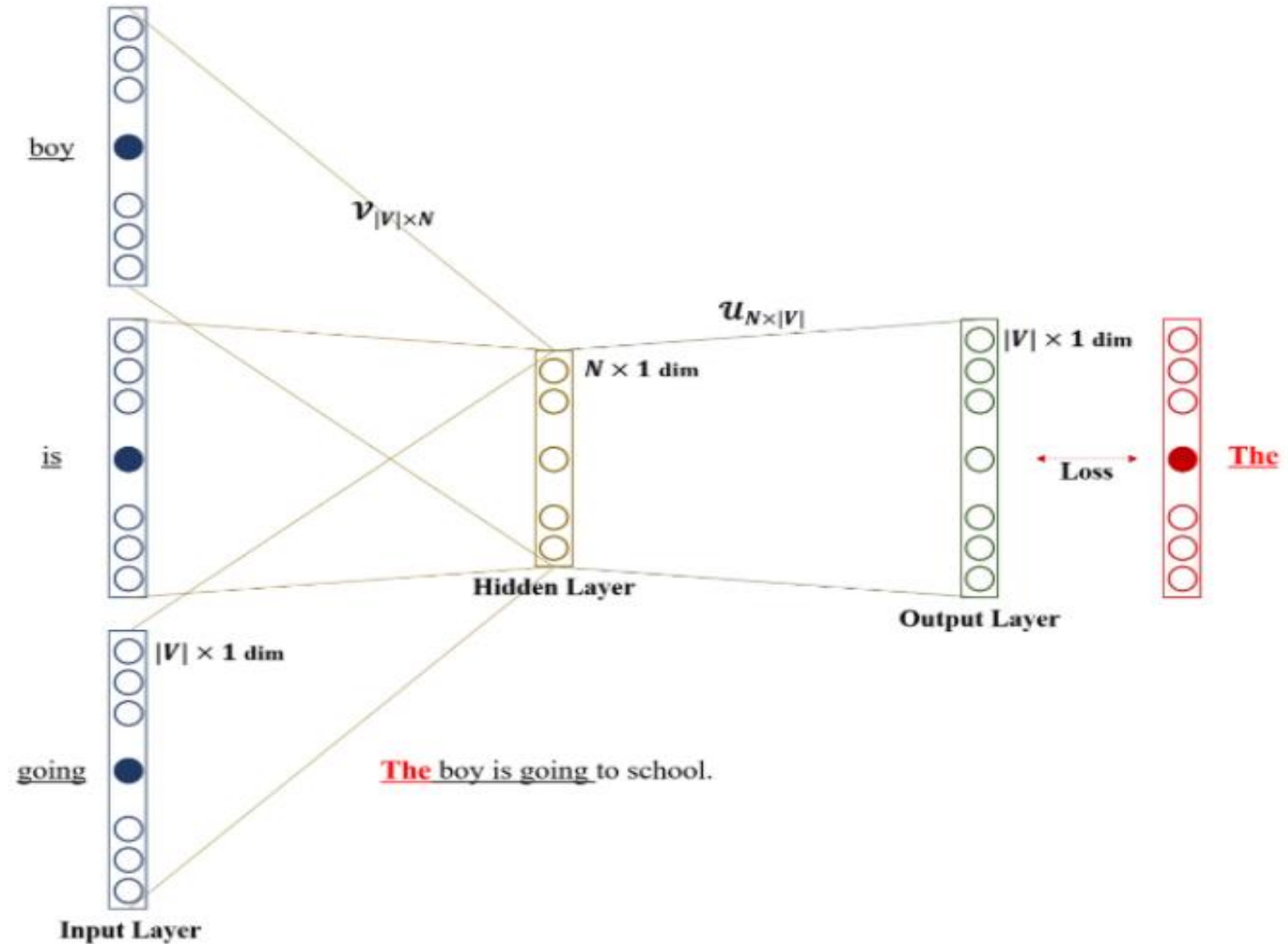
(like, I)
(like, playing)
(like, football)

(playing, I)
(playing, like)
(playing, football)
(playing, with)

..
..
..
..
..
..
..
..
..
..
..

Word2Vec - CBOW

- CBOW architecture



Word2Vec - CBOW

- CBOW Training

1. (Input Layer) **context word** 각각의 one-hot vector를 입력
2. (Input Layer -> Hidden Layer) $v = V^T * \text{input} \rightarrow$ (N-dim vector)
3. (Hidden Layer) averaging v's of Step.2 $\rightarrow v_{\text{hat}}$
4. (Hidden Layer -> Output Layer) $z = U * v_{\text{hat}}$
5. (Output Layer) $y_{\text{hat}} = \text{softmax}(z)$
6. (Loss) error = $y_{\text{hat}} - y(\text{center word})$

$$x^{(c-m)}, \dots, x^{(c-1)}, x^{(c+1)}, \dots, x^{(c+m)} \in \mathbb{R}^{|V|}$$

$$v_{c-m} = \mathcal{V}x^{(c-m)}, \dots, v_{c+m} = \mathcal{V}x^{(c+m)} \in \mathbb{R}^n$$

$$\hat{v} = \frac{v_{c-m} + \dots + v_{c+m}}{2m} \in \mathbb{R}^n$$

$$z = \mathcal{U}\hat{v} \in \mathbb{R}^{|V|}$$

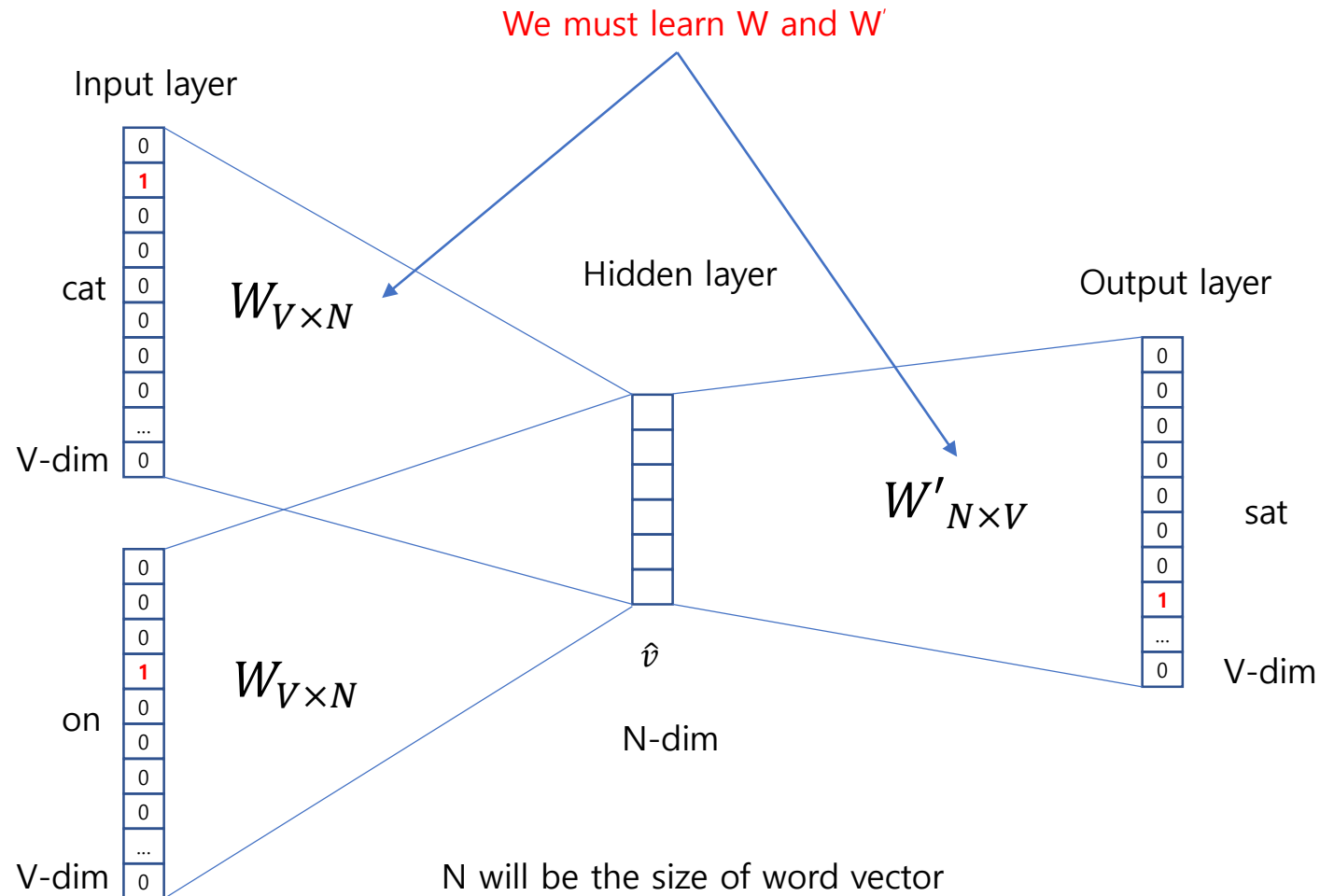
$$\hat{y} = \text{softmax}(z) \in \mathbb{R}^{|V|}$$

$$\mathcal{L}(\hat{y}, y) = - \sum_{j=1}^{|V|} y_j \log \hat{y}_j = -y_i \log \hat{y}_i$$

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \mathcal{L}(\hat{y}, y)$$

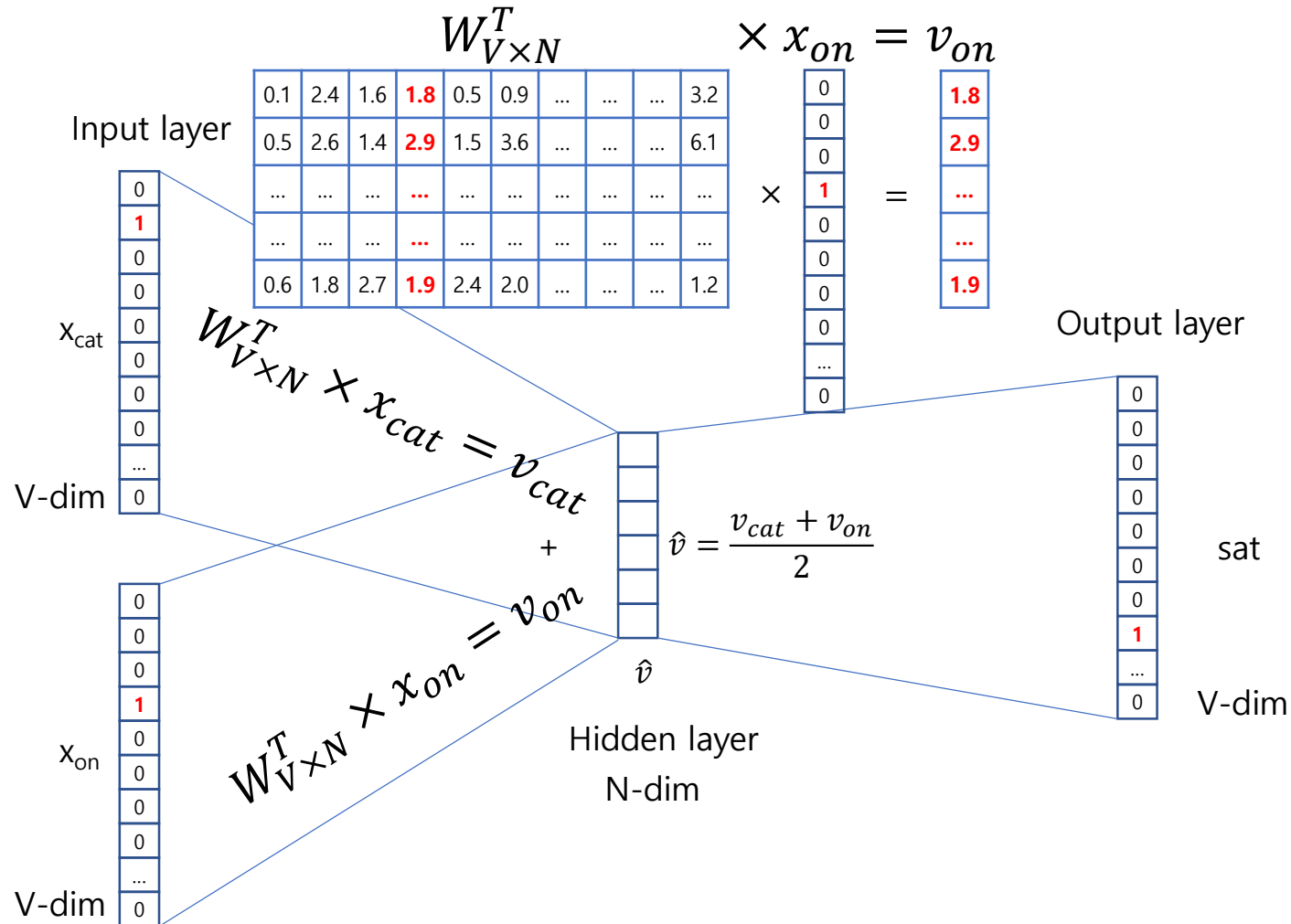
Word2Vec - CBOW

- (ex) "The cat sat on floor" with window_size = 1



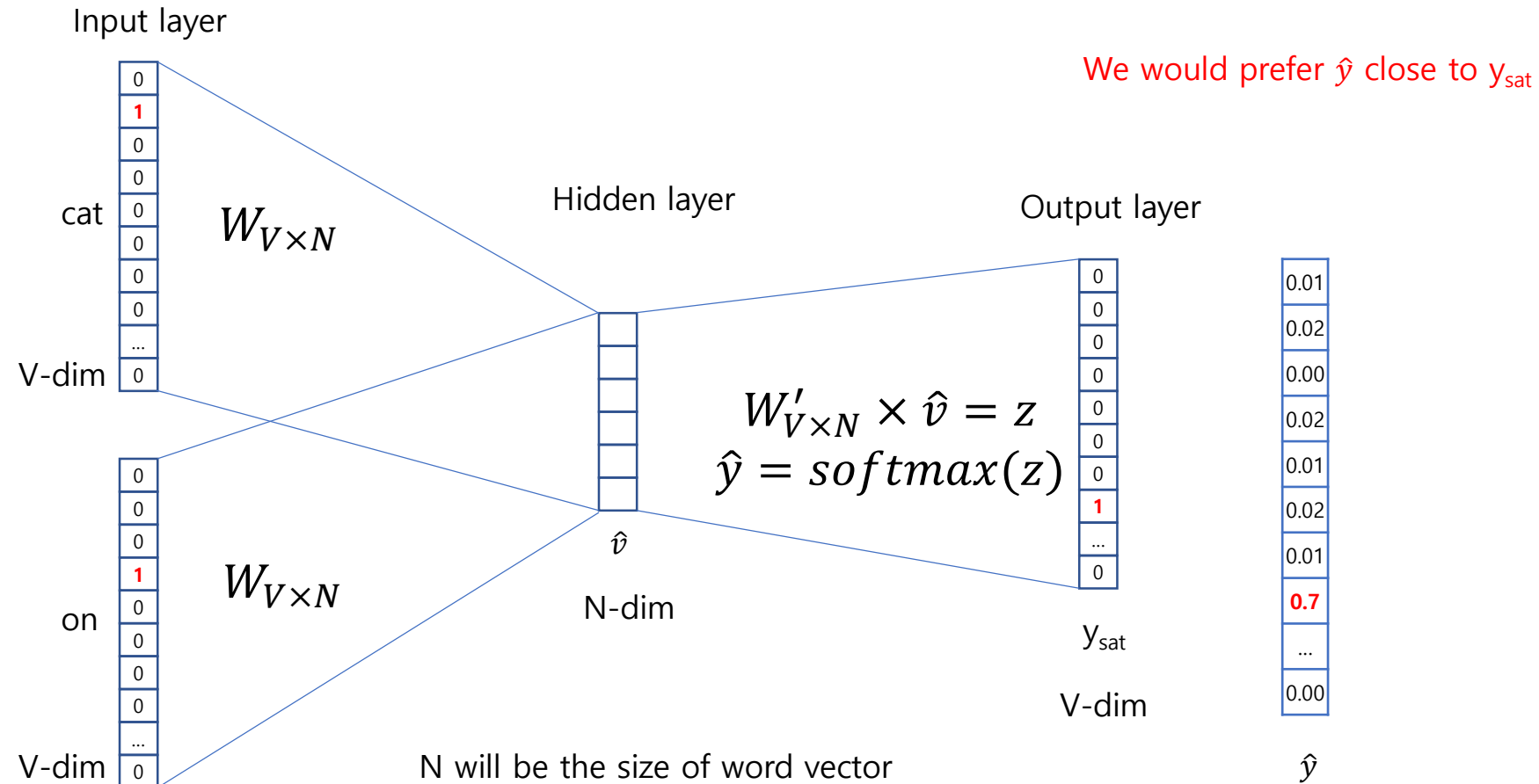
Word2Vec - CBOW

- (ex) "The cat sat on floor" with window_size = 1



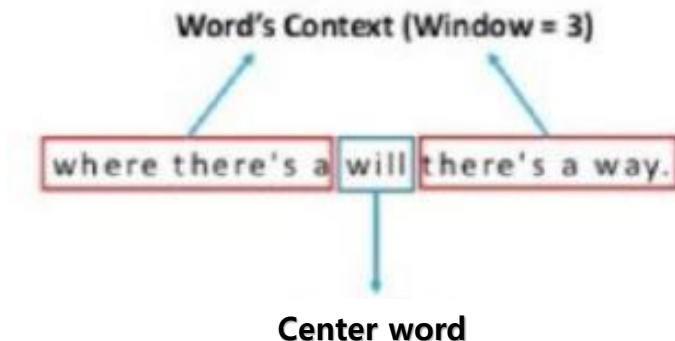
Word2Vec - CBOW

- (ex) "The cat sat on floor" with window_size = 1

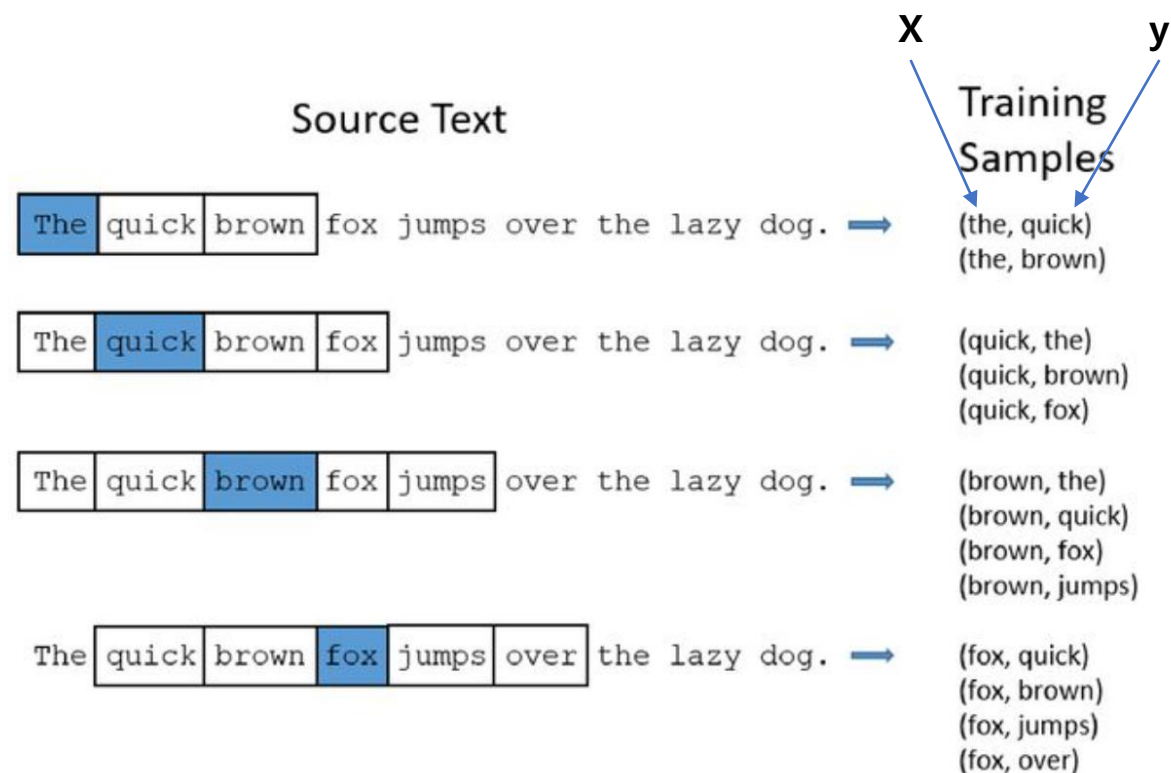


Word2Vec – Skip Gram

- **Predict neighbors** of a center word using Skip-gram model

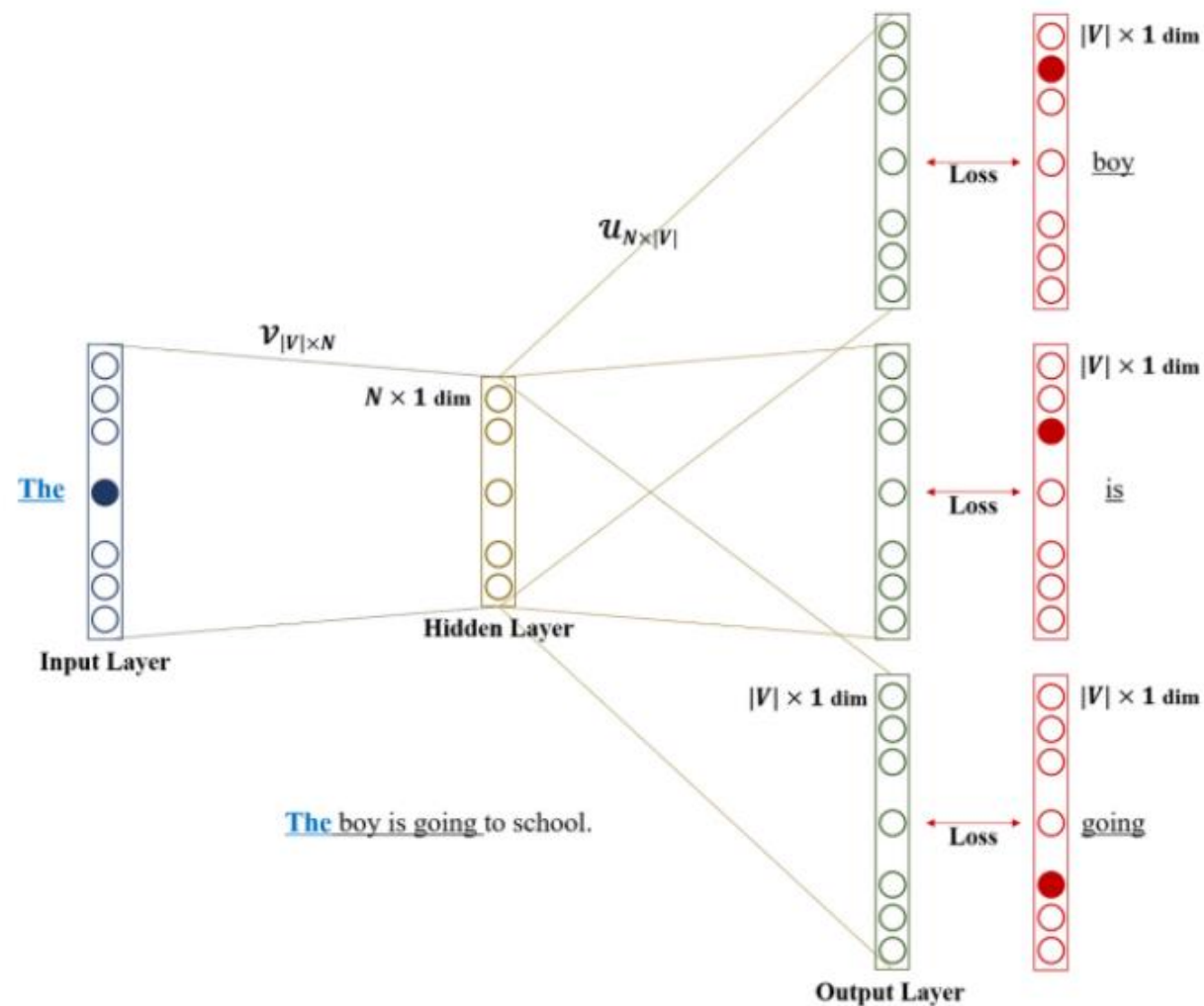


- **Collection of Training samples**
 - with a window of size 2



Word2Vec – Skip Gram

- Skip-Gram architecture



Word2Vec - Skip-Gram

- Skip-Gram Training

1. (Input Layer) **center word** 각각의 one-hot vector를 입력

$$x \in \mathbb{R}^{|V|}$$

2. (Input Layer -> Hidden Layer) $v_c = V^T * \text{input} \rightarrow$ (N-dim vector)

$$v_c = Vx \in \mathbb{R}^n$$

3. (Hidden Layer -> Output Layer) $z = U * v_c$

$$z = Uv_c \in \mathbb{R}^{|V|}$$

4. (Output Layer) $y_{\text{hat}} = \text{softmax}(z)$

$$\hat{y} = \text{softmax}(z) \in \mathbb{R}^{|V|}$$

5. error $e_c = y_{\text{hat}} - y(\text{context word})$

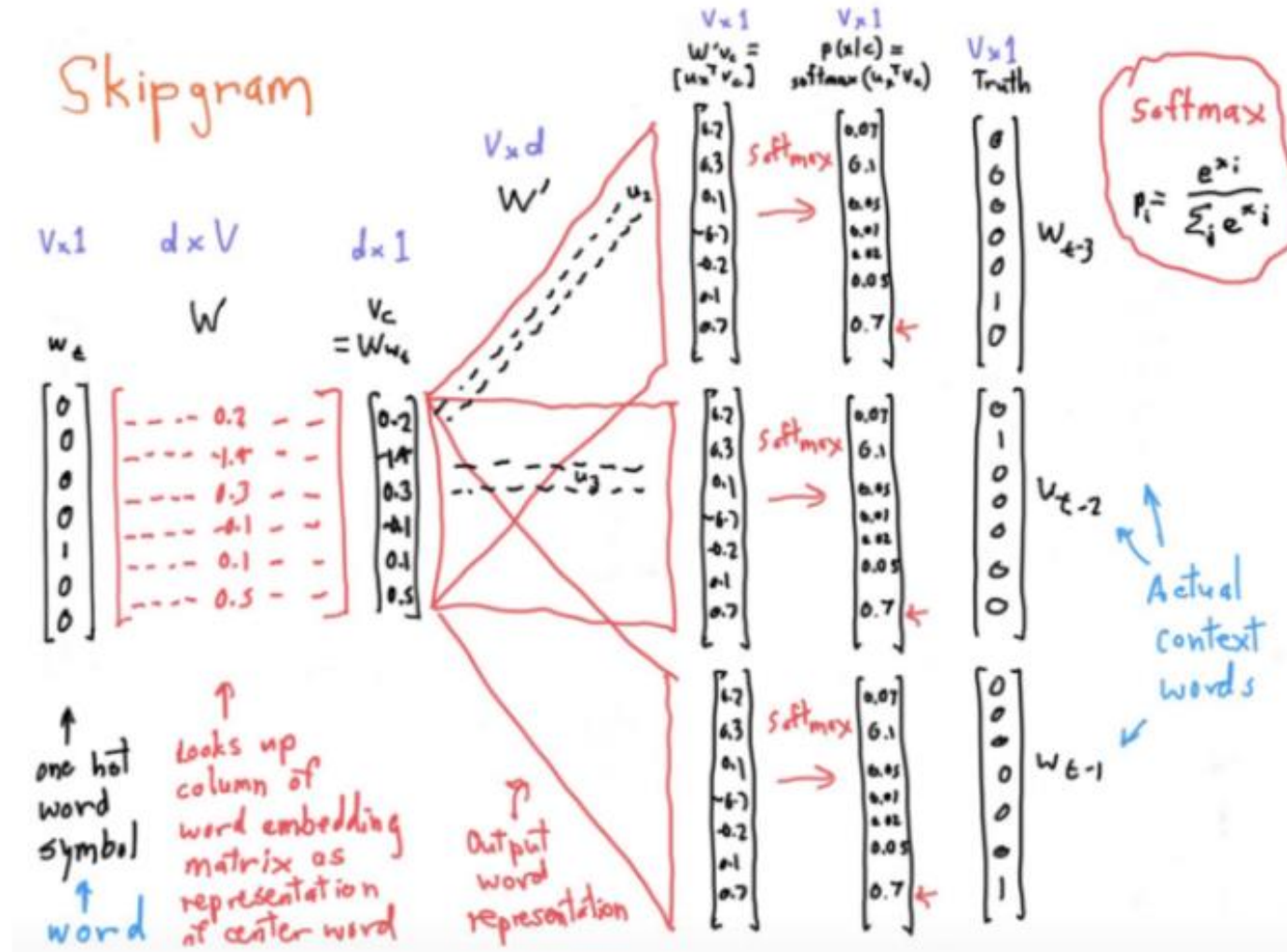
위의 Step 을 모든 context word (주변단어) 에 대해 실행

6. (Loss) prediction error = $\sum_{2m} (e_c)$

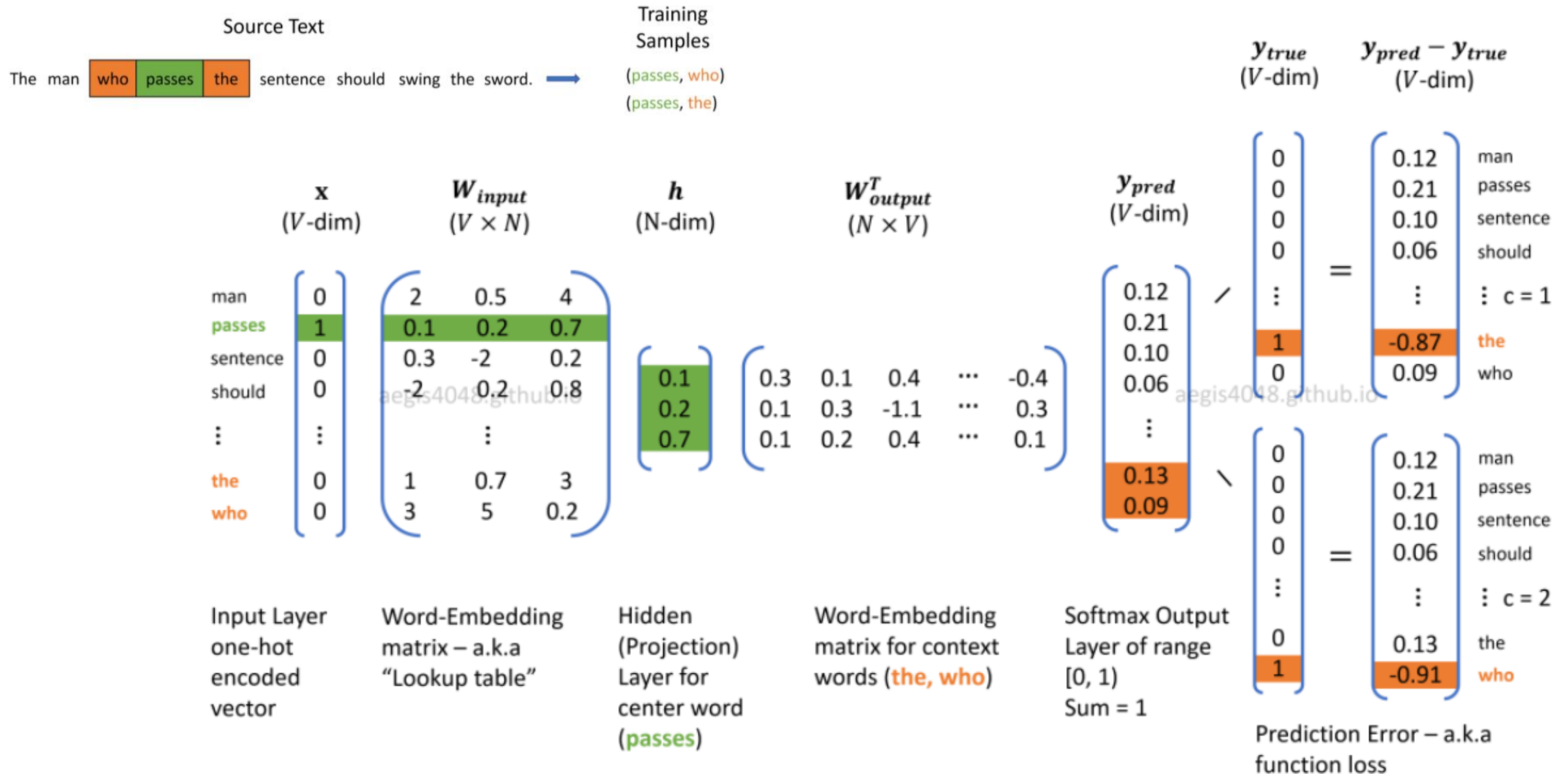
$$\mathcal{L}(\hat{y}, y) = - \sum_{j=0, j \neq m}^{2m} \sum_{k=1}^{|V|} y_k^{(c-j)} \log \hat{y}_k^{(c-j)} = - \sum_{j=0, j \neq m}^{2m} y^{(c-j)} \log \hat{y}^{(c-j)}$$
$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \mathcal{L}(\hat{y}, y)$$

Word2Vec – Skip Gram

- Summary



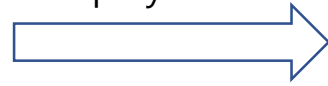
Skip Gram – model example



Skip Gram – Cost function

$$\operatorname{argmax}_{\theta} p(w_1, w_2, \dots, w_C | w_{center}; \theta)$$

to simplify derivation



$$\operatorname{argmax}_{\theta} \log p(w_1, w_2, \dots, w_C | w_{center}; \theta)$$

softmax, naïve assumption



$$\operatorname{argmax}_{\theta} \log \prod_{c=1}^C \frac{\exp(W_{output(c)} \cdot h)}{\sum_{i=1}^V \exp(W_{output(i)} \cdot h)}$$

h : hidden layer word vector

$$p(w_{context} | w_{center}; \theta) = \frac{\exp(W_{output(context)} \cdot h)}{\sum_{i=1}^V \exp(W_{output(i)} \cdot h)}$$

max -> min



$$J(\theta; w^{(t)}) = -\log \prod_{c=1}^C \frac{\exp(W_{output(c)} \cdot h)}{\sum_{i=1}^V \exp(W_{output(i)} \cdot h)}$$

c : index of the context words around the center word (w_t)
 C : window size

$$= -\sum_{c=1}^C \log \frac{\exp(W_{output(c)} \cdot h)}{\sum_{i=1}^V \exp(W_{output(i)} \cdot h)} = -\sum_{c=1}^C (W_{output(c)} \cdot h) + C \cdot \log \sum_{i=1}^V \exp(W_{output(i)} \cdot h)$$

Skip Gram – Cost function

$$W_{input}^{(new)} = W_{input}^{(old)} - \eta \cdot \frac{\partial J}{\partial W_{input}}$$

$$\frac{\partial J}{\partial W_{input}} = x \cdot (W_{output}^T \sum_{c=1}^C e_c)$$

$$W_{output}^{(new)} = W_{output}^{(old)} - \eta \cdot \frac{\partial J}{\partial W_{output}}$$

$$\frac{\partial J}{\partial W_{output}} = h \cdot \sum_{c=1}^C e_c$$



$$W_{input}^{(new)} = W_{input}^{(old)} - \eta \cdot x \cdot (W_{output}^T \sum_{c=1}^C e_c)$$

$$W_{output}^{(new)} = W_{output}^{(old)} - \eta \cdot h \cdot \sum_{c=1}^C e_c$$

$\sum_{c=1}^C e_c =$

$c = 1$	$c = 2$	
0.12	0.12	man
0.21	0.21	passes
0.10	0.10	sentence
0.06	0.06	should
\vdots	\vdots	\vdots
-0.87	0.13	the
0.09	-0.91	who

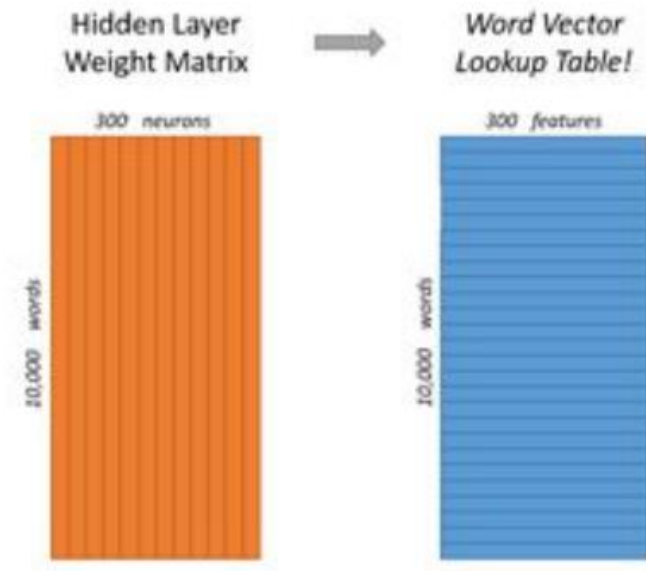
prediction error

Word Embedding

- **Build a vocabulary of words from training documents**
 - (e.g.) a vocabulary of 10,000 unique words
- **Represent an input word as a one-hot vector**
 - This vector will have 10,000 components (one for every word in our vocabulary)
 - This vector will have "1" in the position corresponding to the word, say "ants", and 0s in all of the other positions.
- **No activation function for hidden layer neurons, but output neurons use softmax.**
- When *training* the network with word pairs, **the input is a one-hot vector and output is also a one-hot vector** representing the output word.
- When evaluating the network on an input word, the output vector will actually be a probability distribution (i.e., a bunch of floating point values, **not** a one-hot vector).

Word Embedding

- Hidden layer is represented by a weight matrix with 10,000 rows (one for every word in our vocabulary) and 300 columns (one for every hidden neuron)
 - 300 features
- Rows of this weight matrix are the word vectors!



Word Embedding

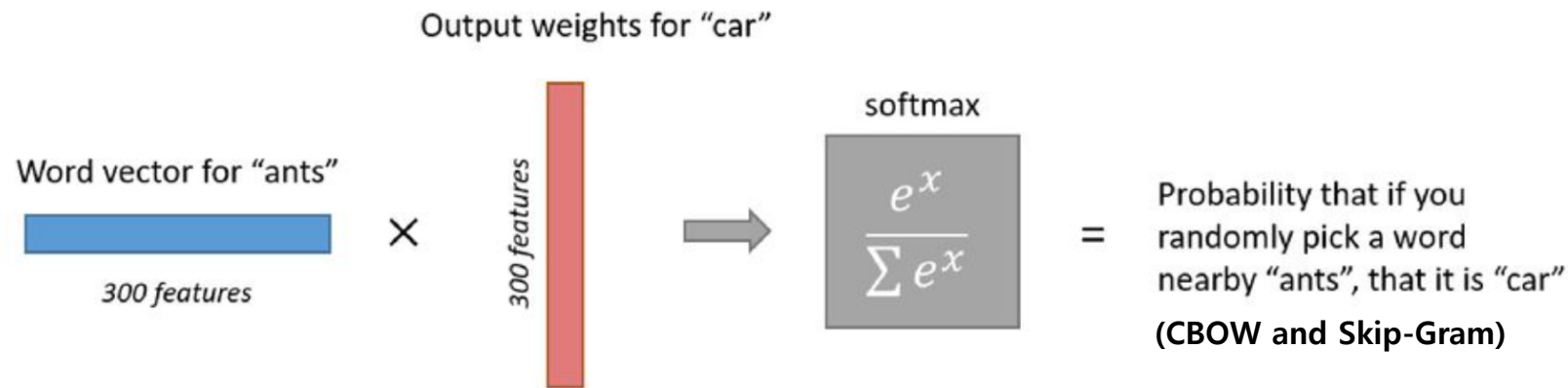
- If we multiply a 1x10,000 one-hot vector by a 10,000x300 matrix, it will effectively just *select* the matrix row corresponding to the “1”

$$[0 \quad 0 \quad 0 \quad 1 \quad 0] \times \begin{bmatrix} 17 & 24 & 1 \\ 23 & 5 & 7 \\ 4 & 6 & 13 \\ 10 & 12 & 19 \\ 11 & 18 & 25 \end{bmatrix} = [10 \quad 12 \quad 19]$$

- Hidden layer is really just operating as a lookup table !!
- The output of the hidden layer is just the “word vector” for the input word.

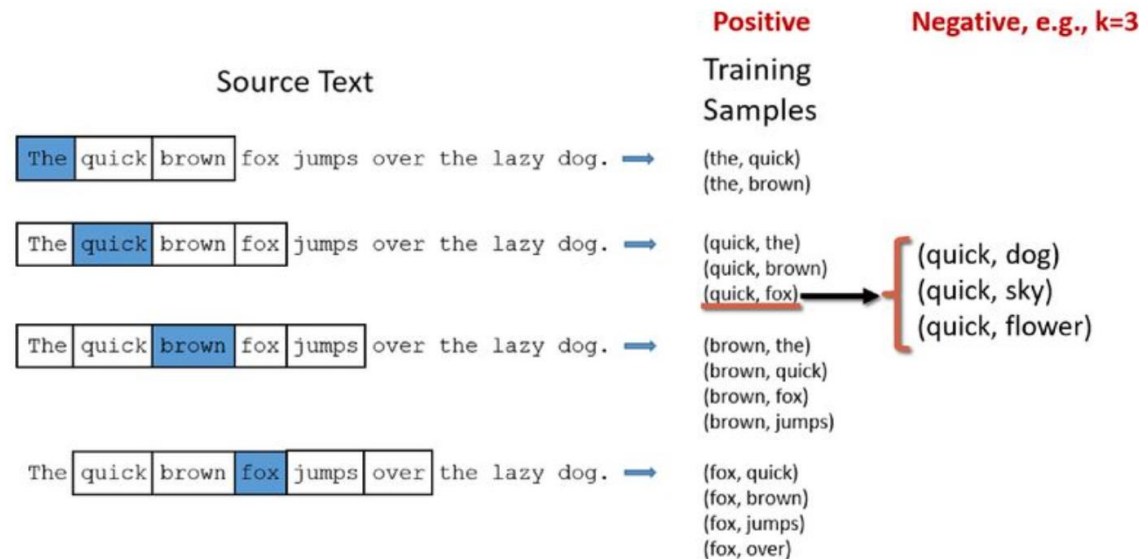
Word Embedding

- The output layer is a **softmax regression classifier**
- Each output neuron has a weight vector which it multiplies against the word vector from the hidden layer, then it applies the function $\exp(x)$ to the result. Finally, in order to get the outputs to sum up to 1, we divide this result by the sum of the results from *all* 10,000 output nodes.



Negative Sampling

- Softmax function requires too much computation for its denominator (say, summation of 10,000 terms)
- When training, **negative sampling** considers only a small number of negative words (let's say 5) including the positive word for this summation (normally choose more than 20 based on some probability distribution)
 - "negative" word is one for which network outputs "0" and "positive" word is one for "1"
 - 즉, 파라미터 조정할 때 전체 단어 집합이 아닌 일부 단어만 조정 (주변 단어(positive)들은 모두 조정하지만 그렇지 않은 단어들(negative)은 일부만 조정)



Negative Sampling

- In the Skip-gram case, we build a **new objective function** that tries to maximize both $P(D=1|w,c,\theta)$ and $P(D=0|w,c,\theta)$. We model $P(D=1|w,c,\theta)$ and $P(D=0|w,c,\theta)$ as follows:

$$P(D=1|w,c,\theta) = \sigma(u_c^T v_w) = \frac{1}{1 + e^{-u_c^T v_w}} \quad \leftarrow \text{Probability that the word pair (w,c) is in corpus}$$

$$P(D=0|w,c,\theta) = 1 - P(D=1|w,c,\theta) = 1 - \frac{1}{1 + e^{-u_c^T v_w}} = \frac{1}{1 + e^{u_c^T v_w}} = \sigma(-u_c^T v_w) \quad \leftarrow \text{Probability that the word pair (w,c) is not in corpus}$$

- We take a simple maximum likelihood approach of these two probabilities. (Here we take θ to be the parameters of the model, and in our case it is V and U .)

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

Θ : matrix U, V

Negative Sampling

- Unnatural sentences that should get a low probability of ever occurring. We can generate \tilde{D} on the fly by randomly sampling this negative from the word bank. [Our new objective function](#) would then be:

$$J = - \sum_{(w,c) \in D} \log \frac{1}{1 + \exp(-u_w^T v_c)} - \sum_{(w,c) \in \tilde{D}} \log \left(\frac{1}{1 + \exp(u_w^T v_c)} \right) = \log \sigma(u_{c-m+j}^T \cdot v_c) + \sum_{k=1}^K \log \sigma(-\tilde{u}_k^T \cdot v_c)$$

D: positive sample space

\tilde{D} : negative sample space

- Original Skip-gram

$$\begin{aligned} \text{minimize } J &= -\log P(w_{c-m}, \dots, w_{c-1}, w_{c+1}, \dots, w_{c+m} | w_c) \\ &= -\log \prod_{j=0, j \neq m}^{2m} P(w_{c-m+j} | w_c) \\ &= -\log \prod_{j=0, j \neq m}^{2m} P(u_{c-m+j} | v_c) \\ &= -\log \prod_{j=0, j \neq m}^{2m} \frac{\exp(u_{c-m+j}^T v_c)}{\sum_{k=1}^{|V|} \exp(u_k^T v_c)} \\ &= - \sum_{j=0, j \neq m}^{2m} u_{c-m+j}^T v_c + 2m \log \sum_{k=1}^{|V|} \exp(u_k^T v_c) \end{aligned}$$

- Skip-gram with Negative Sampling (SGNS)

$$\sum_{j=0, j \neq m}^{2m} \log(\sigma(U_{c-m+j} \cdot v_c)) + \sum_{j=1}^{\tilde{K}} \log(\sigma(-\tilde{U}_j \cdot v_c))$$

Only K samples

$$- \sum_{j=0, j \neq m}^{2m} U_{c-m+j} \cdot v_c + 2m \log \sum_{j=1}^{|V|} e^{U_j \cdot v_c}$$

Word Analogies

- Interesting result

Test for linear relationships, examined by Mikolov et al. (2014)

a:b :: c:?



$$d = \arg \max_x \frac{(w_b - w_a + w_c)^T w_x}{\|w_b - w_a + w_c\|}$$

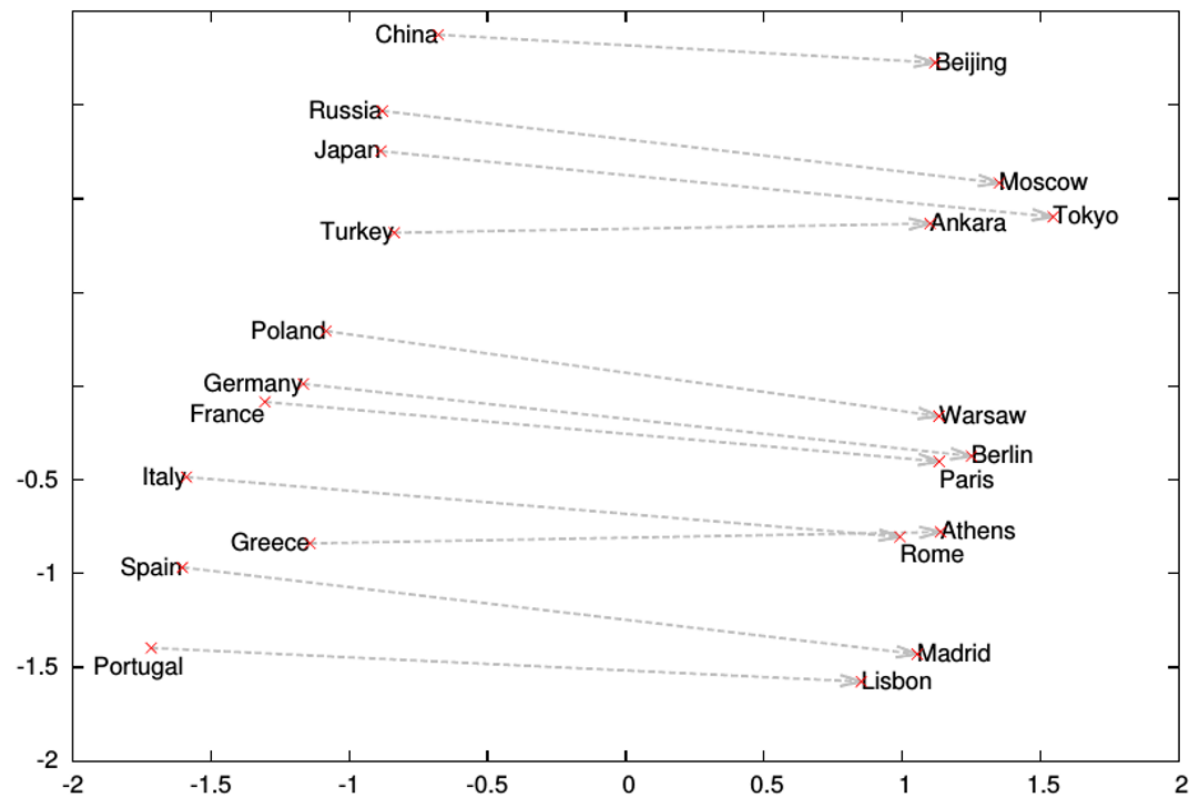
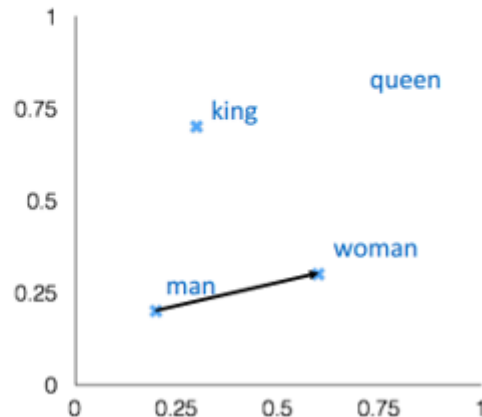
man:woman :: king:?

+ king [0.30 0.70]

- man [0.20 0.20]

+ woman [0.60 0.30]

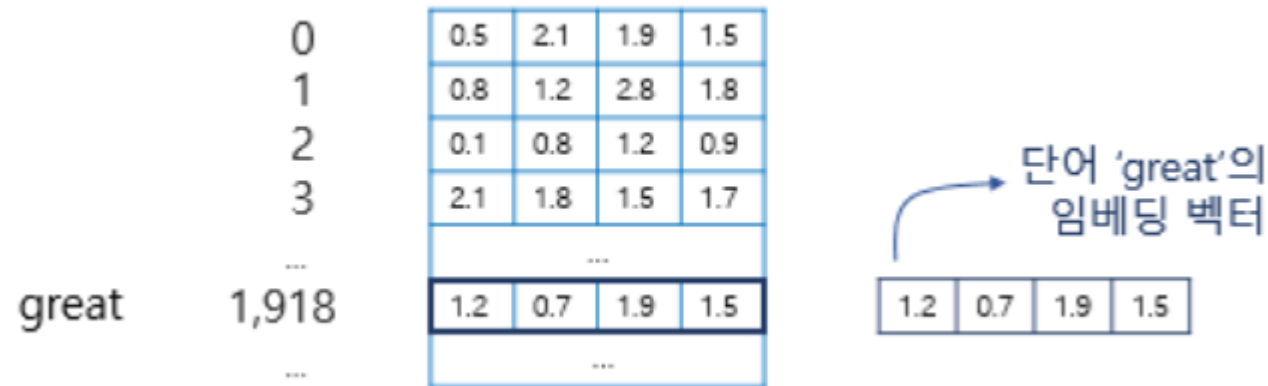
queen [0.70 0.80]



Keras Embedding() Layer

- A word embedding can be **learned as part of a deep learning model**. This can be a slower approach, but tailors the model to a **specific training dataset**.
- Embedding layer is just a **Look-up Table**.
 - All the words should be encoded as integers to be used as input.

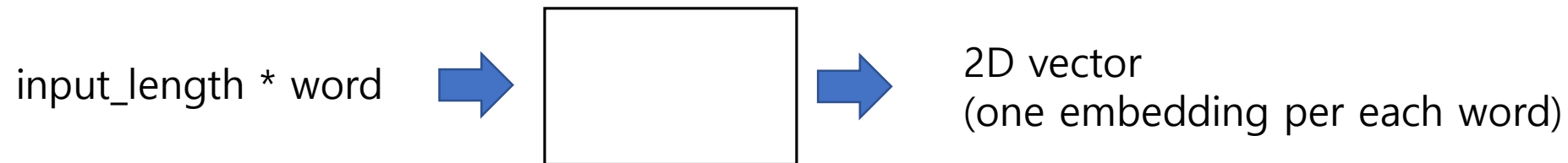
Word → Integer → lookup Table → Embedding vector



훈련 과정에서 학습된다.

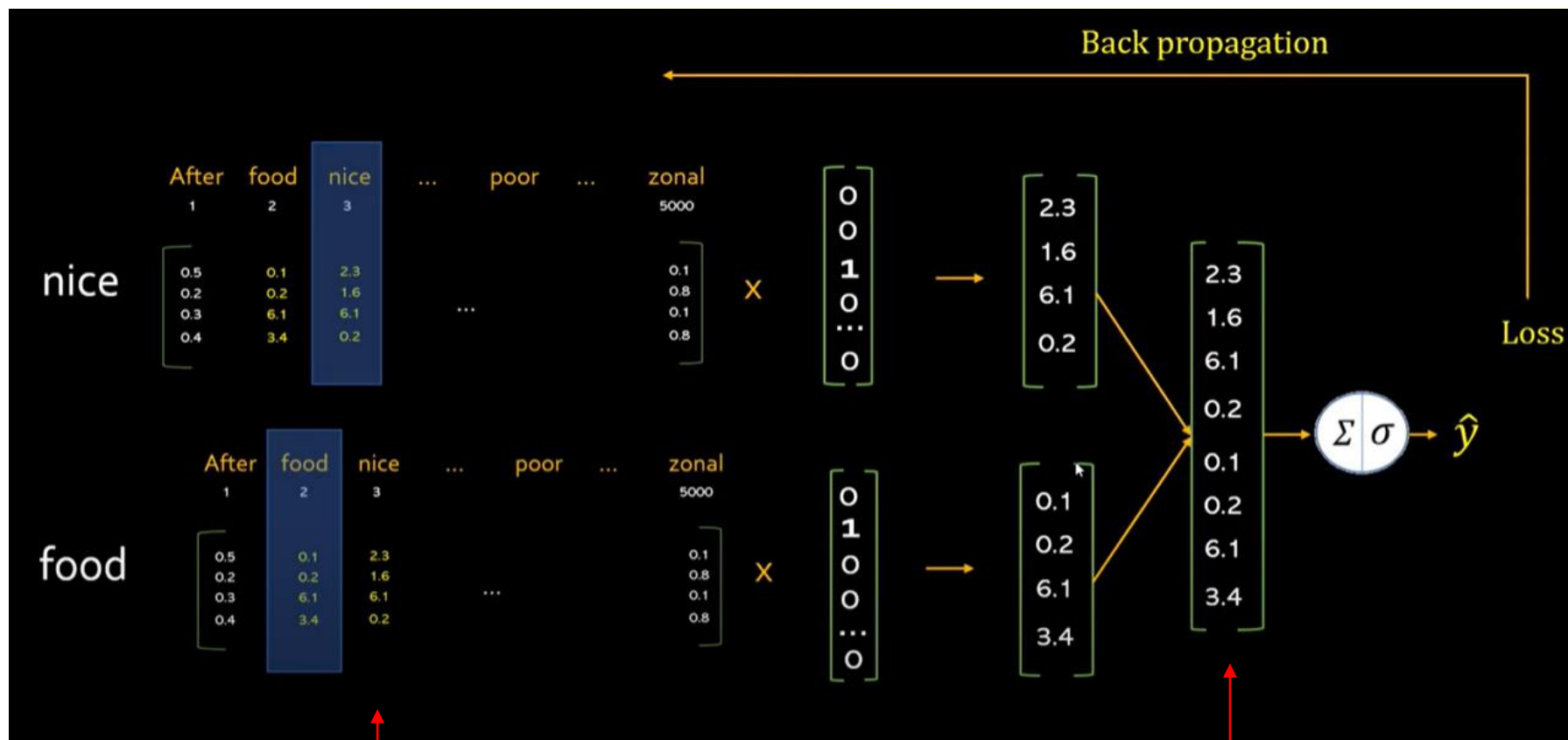
Keras Embedding() Layer

- **Flexible layer**
 - Can be used alone, or can be used as part of deep learning model
 - Can also be used to load a pre-trained word embedding model
- Embedding layer is defined as **the first layer** of a network.
- To be connected directly to Dense layer, the 2D output must be flattened to 1D.
- **3 arguments**
 - `input_dim`: size of the vocabulary
 - `output_dim`: size of the vector space
 - `input_length`: length of input sequences (e.g. if all of your input documents are comprised of 1000 words, this would be 1000.)



(*) The result of Embedding is just a permutation of inner trainable weights.
Where the permutation is denoted by "indices" in your input array. (see the next slide)

Inside of Keras Embedding() Layer

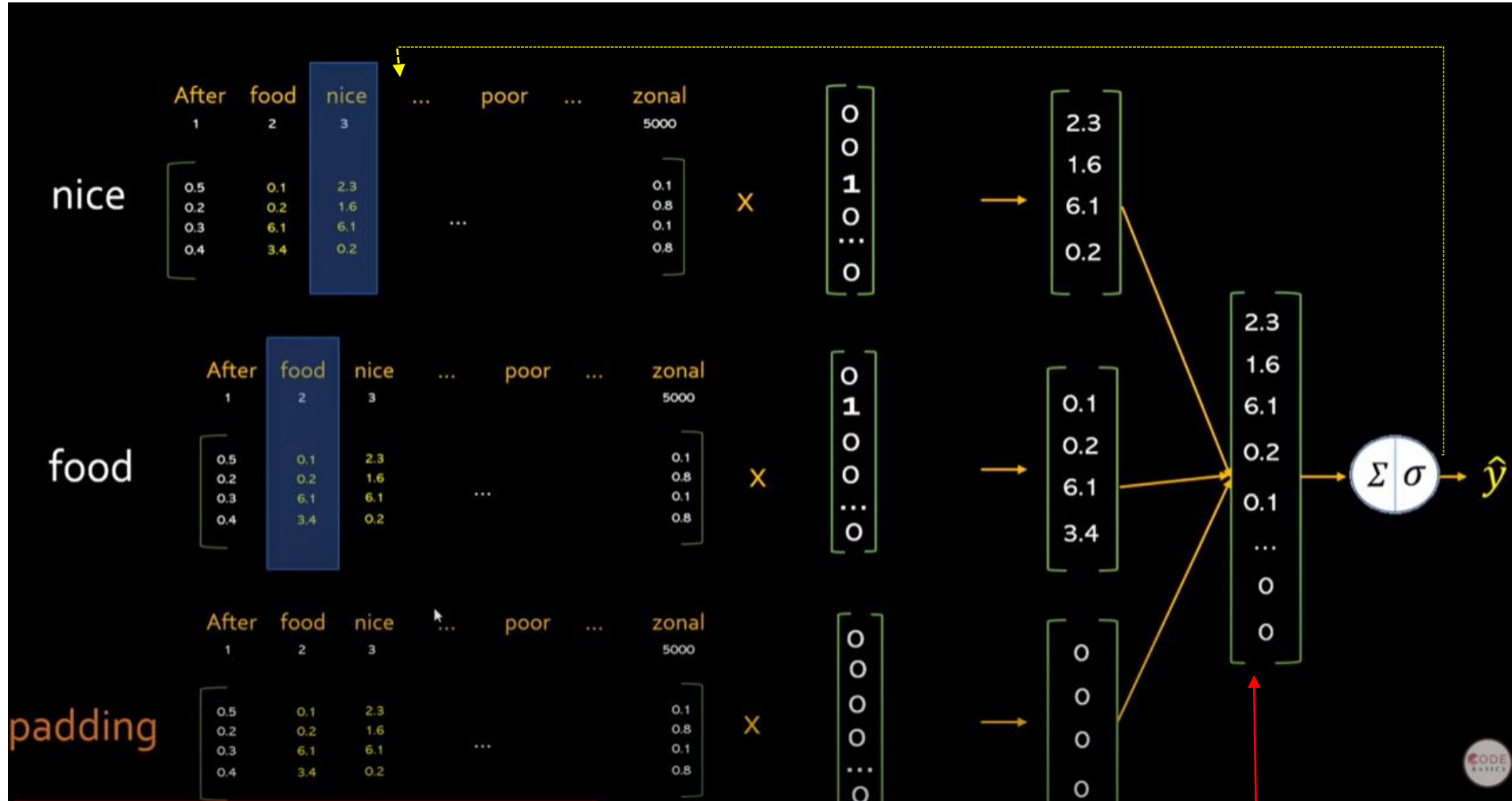


Values of the matrix are initially randomly selected, and will keep changing during the training stage

number of neurons is fixed

<https://www.youtube.com/watch?v=Fuw0wv3X-0o>

Inside of Keras Embedding() Layer



Need padding to fix the network

Inside of Keras Embedding() Layer

- The matrix after training



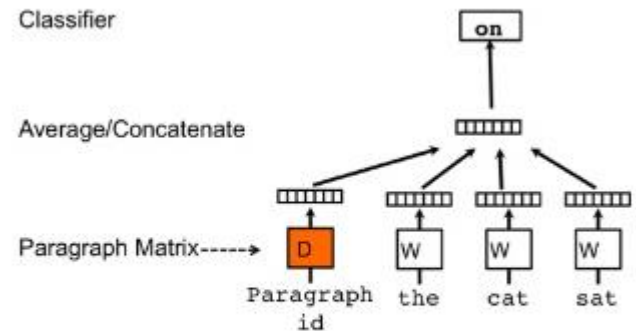
similar words are embedded similarly

Doc2Vec

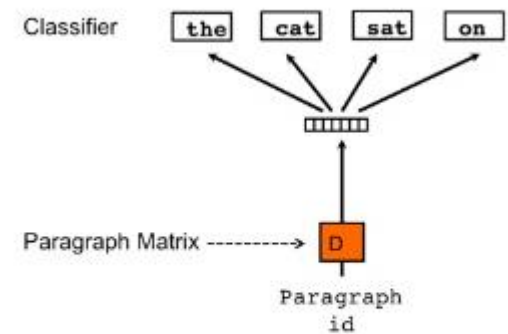
- **Embedding the paragraphs (documents) to vectors**
- **Extension of Word2Vec**
- **Basic Idea**
 - act as if a document has **another floating word-like vector**, which contributes to all training predictions, and is updated like other word-vectors, but we will call it a **doc-vector**.
- **Two implementations**
 - Paragraph Vector – Distributed Memory (PV-DM)
 - Paragraph Vector – Distributed bag of Words (PV-DBOW)

Doc2Vec

- **Example: “The cat sat on the mat.” in paragraph_1**
 - (Training data)
 - window size $k = 3$
 - [paragraph_1, the cat, sat] – on
 - [paragraph_1, cat, sat, on] – the
 - [paragraph_1, sat, on, the] – mat
- **PY-D2V and PV-DBOW**



Paragraph vector with distributed memory (PV-DM)



Distributed bag of words version (PVD BOW)

Topic Modeling

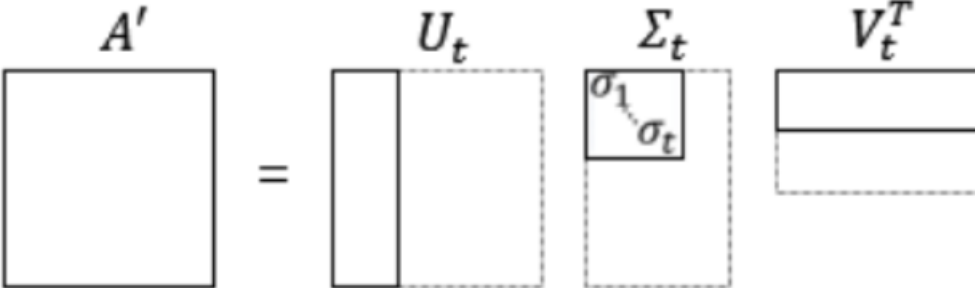
- **Topic Modeling:** a technique used to discover [underlying topics or themes](#) within a collection of documents.
 - It assumes that documents are mixtures of different topics, and topics are distributions over words.
 - The goal is to identify these topics and their distributions in a way that helps uncover the main themes present in the text data.
- **Latent Semantic Analysis (LSA):** LSA is a technique that uses SVD to perform topic modeling. It operates on a term-document matrix (TDM). By applying SVD to this matrix, LSA identifies latent (hidden) topics by capturing patterns of co-occurrence in the data.
- **Latent Dirichlet Allocation (LDA):** a generative probabilistic model used in natural language processing and machine learning.

Topic Modeling

- **LSA (Latent Semantic Analysis)**

- LDA (토픽 모델링) 에 아이디어 제공한 알고리즘
- DTM(document-term matrix) 나 TfIdf matrix 에 Truncated SVD 이용
- DTM을 차원 축소 하여 축소 차원에서 근접 단어들을 토픽으로 묶는다.

Truncated SVD

$$A' = U_t \Sigma_t V_t^T$$


- **LDA (Latent Dirichlet Allocation)**

- 각 단어가 특정 토픽에 존재할 확률과 문서에 특정 토픽이 존재할 확률을 결합확률로 추정하여 토픽을 추출하여 할당

LSA (Latent Semantic Analysis)

- Matrix decomposition

$$\begin{array}{c} \left[\begin{array}{c} M \times K \end{array} \right] \\ \text{Topic Assignment} \end{array} \times \begin{array}{c} \left[\begin{array}{c} K \times V \\ \text{Topics} \end{array} \right] \approx \begin{array}{c} \left[\begin{array}{c} M \times V \\ \text{Dataset} \end{array} \right] \end{array}$$

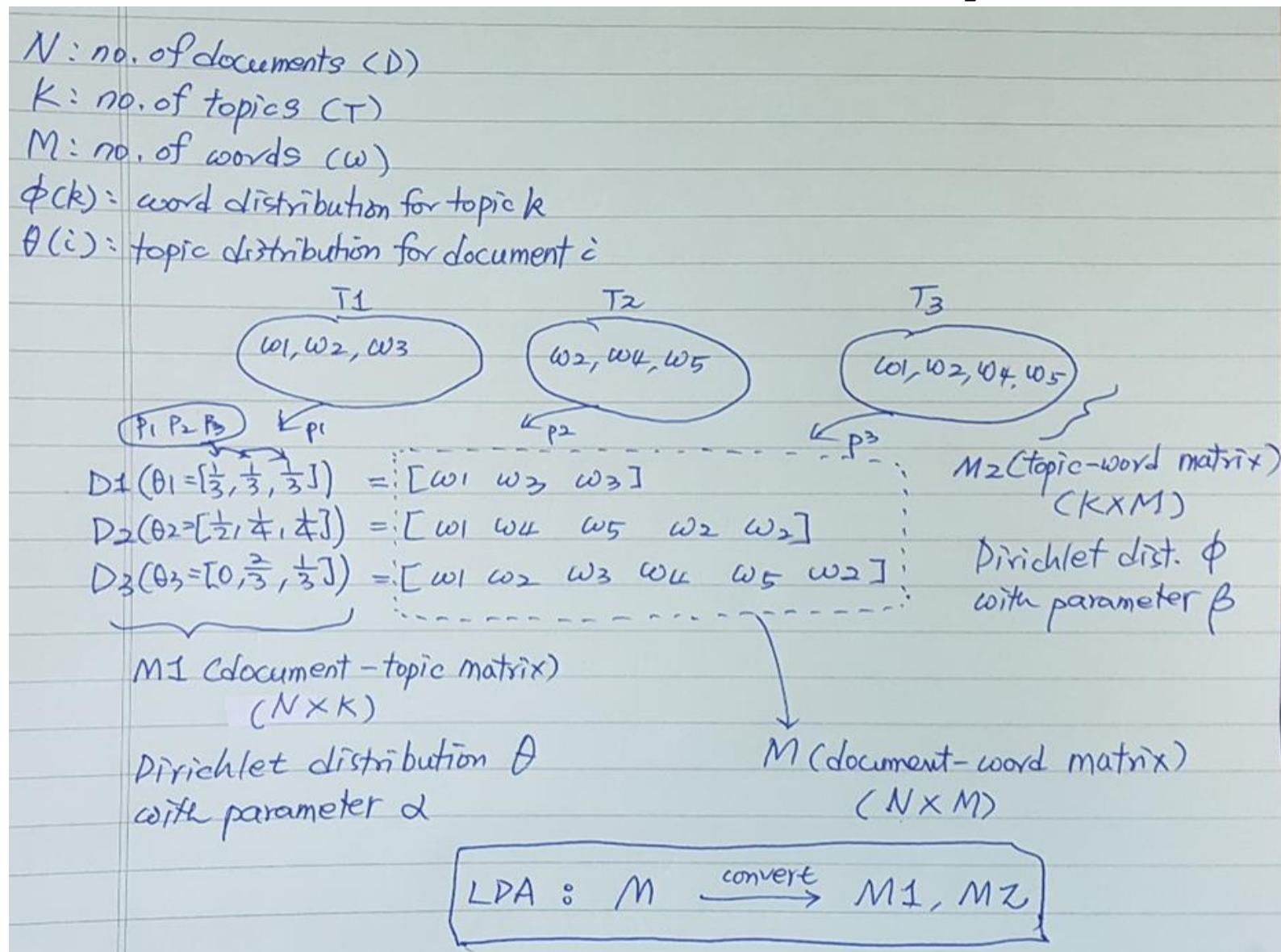
K: Number of topics
M: Number of documents
V: number of words

- $M = U * \Sigma * VT$

- Each row of U can be seen as a vector that represents the document in terms of the underlying topics (This means that each document is associated with a distribution of topics, where some topics might be more prominent than others for that particular document.)
- Each column of VT can be interpreted as a vector that represents a term's contribution to the various topics. (This means that each term is associated with a distribution of topics, indicating which topics it is most relevant to.)
- The singular values in Σ don't have a direct semantic interpretation in terms of topics or documents. (The primary role of Σ in LSA is in the reduction of dimensions.)
- Only the **relative magnitudes** of values within the matrices are important (not the absolute values). -> Larger values represent stronger associations.

LDA (Latent Dirichlet Allocation)

- Assumption:
 - The documents (D) are generated by the procedure shown in the memo.
- LDA
 - Backtracking the procedure
 - In LDA, you estimate the parameters α and β from the observed data (i.e. the documents) using techniques like Gibbs sampling.



LDA (Latent Dirichlet Allocation)

- Dirichlet Distribution (Probabilistic k-simplex)

$$f(x_1, \dots, x_k; \alpha_1, \dots, \alpha_k) = \frac{1}{B(\alpha)} \prod_{i=1}^k x_i^{\alpha_i - 1}$$

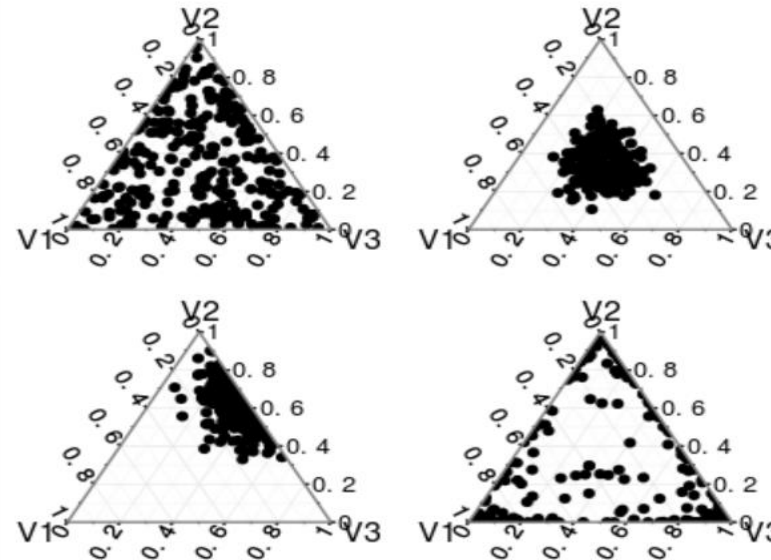
$$B(\alpha) = \frac{\prod_{i=1}^k \Gamma(\alpha_i)}{\Gamma(\sum_{i=1}^k \alpha_i)} \quad (\Gamma \text{ is gamma function})$$

continuous random variables

$$x_k, \sum_k x_k = 1, \forall x_k \geq 0$$

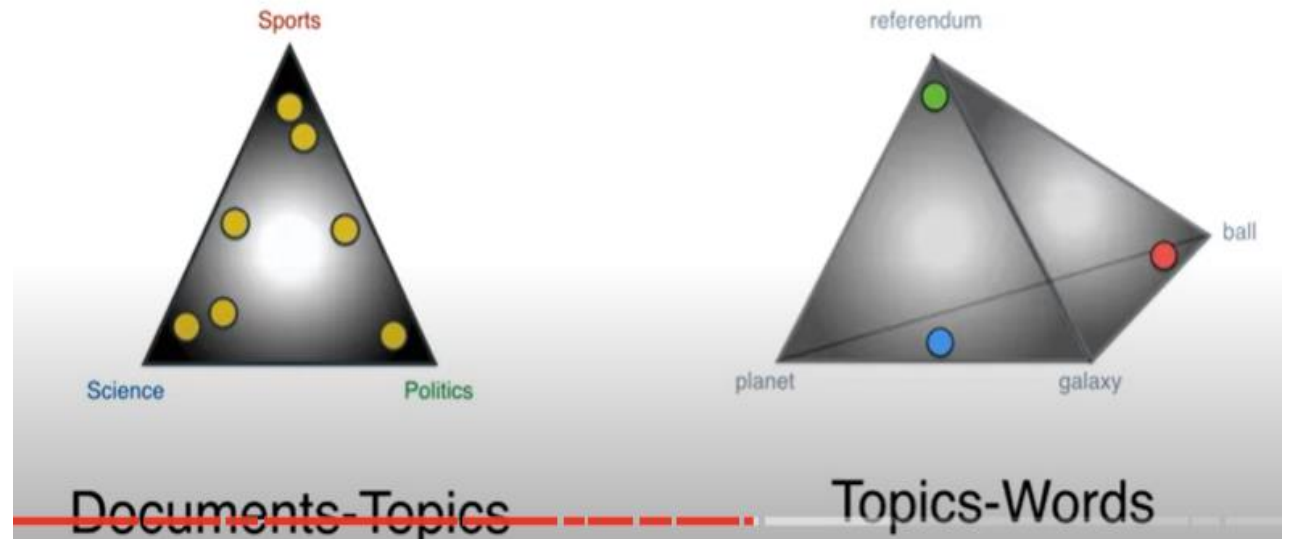
- Example (k=3)

- (a) $\alpha_1 = \alpha_2 = \alpha_3 = 1$
- (b) $\alpha_1 = \alpha_2 = \alpha_3 = 10$
- (c) $\alpha_1 = 1, \alpha_2 = 10, \alpha_3 = 5$
- (d) $\alpha_1 = \alpha_2 = \alpha_3 = 0.2$



LDA

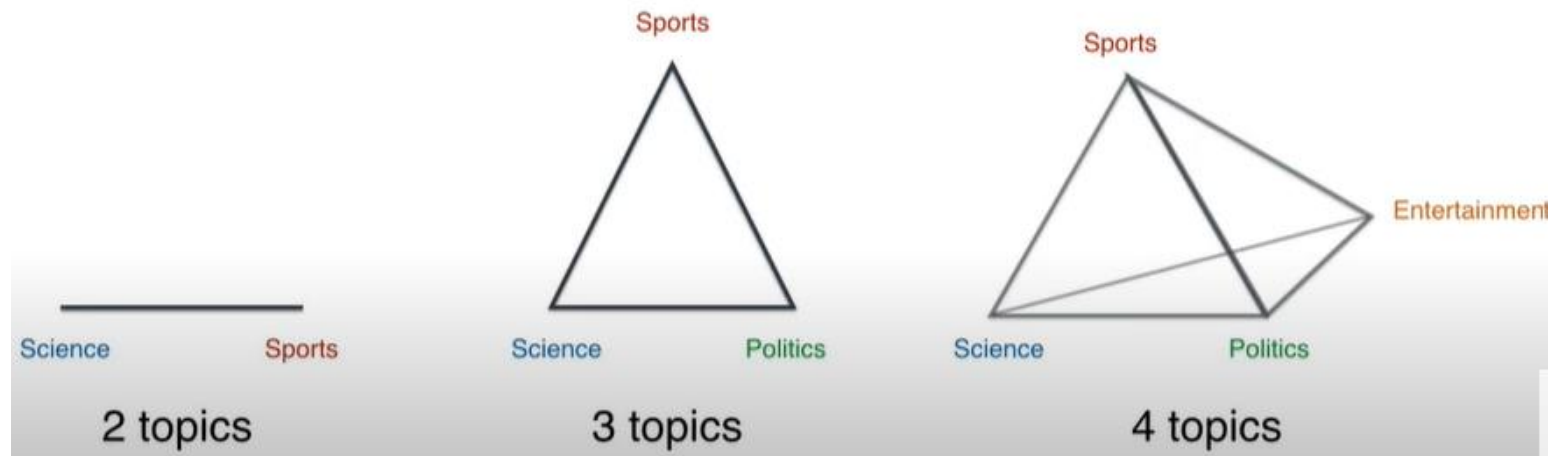
- **Goal:**
 - Assign an appropriate Topic to each Word
- **Properties:**
 - Each document is as monochromatic as possible
 - Each word is as monochromatic as possible
- **Probabilistic Topic model**
 - Each document is a probability distribution over topics
 - Each topic is a probability distribution over words
- **Two Dirichlet distributions**
 - (ex) 3 Topics (Red, Green, Blue)
 - 7 documents with 4 words (referendum, ball, galaxy, planet)



LDA (Latent Dirichlet Allocation)

- Dirichlet Distribution

(from https://www.youtube.com/watch?v=BaM1uiCpj_E)



LDA

- Gibbs sampling (example)



1. Randomly assign Topics to each word
2. By assuming that all other words are properly assigned, reassign (probabilistically) all the words (one by one) based on:
 - 단어가 속한 그 document 에 topic 들이 어떻게 분포되어 있는지
 - 단어가 전체 documents 에서 각 topic 들에 어떻게 assign 되어 있는지
3. Classify each document for the topics based on majority.

LDA

- A generative probabilistic model, and while it is often used in the context of clustering and topic modeling.
- LDA is a model for discovering topics within a collection of documents.
- It assumes that documents are mixtures of topics, and words within documents are generated from these topics.
- How LDA differs from traditional clustering models:
 - Topic Modeling: its primary goal is to uncover the underlying topics present in the document collection and their distribution across documents.
 - Probabilistic Generative Model: it provides a probabilistic framework for generating documents. (it simulates how documents are created based on topic proportions and word distributions)
 - Soft Assignment: each document has a probability distribution over topics, and each word has a probability distribution over topics as well. (not hard-assigned to a single cluster or topic)
 - Bayesian Inference: employs Bayesian inference techniques, such as Gibbs sampling, to estimate the model's parameters, including the topic distributions and word-topic assignments

LDA: two implementations

	sklearn	gensim
Focus	a simpler implementation (primary focus is on dimensionality reduction techniques like PCA)	more specialized for natural language processing tasks, including topic modeling (core topic modeling algorithm)
Corpora and Dictionaries	requires a term-document matrix (TF-IDF matrix) as input. It doesn't include utilities for creating corpora or dictionaries.	Provides tools for creating and working with corpora and dictionaries. It uses its own internal representation called corpus and dictionary objects.
Input Data Format	term-document matrix (TF-IDF matrix)	Expects a corpus object, which is a collection of documents represented as lists of (word ID, word frequency) tuples. It also requires a dictionary mapping word IDs to words.
Model Parameters	fewer options for configuring LDA, which makes it easier for beginners	more flexible (more control over model parameters and hyperparameters)
Additional Functionality	focuses on dimensionality reduction and simpler	Offers various additional functionalities for topic modeling, such as computing topic coherence, visualization using pyLDAvis, and more.
Training	Batch approach Use standard interface: fit_transform()	On-line approach (model is trained iteratively using smaller batches of documents) When you initialize the LdaModel, it starts with random topic assignments for words. Then, you use the LdaModel's train() method to iteratively update the topic and document distributions until convergence.