man

woman

king

queen

Male-Female

walked

swam

walking

swimming

Verb tense

Spain

Italy — Madrid

Germany — Rome

— Berlin

Turkey

— Ankara

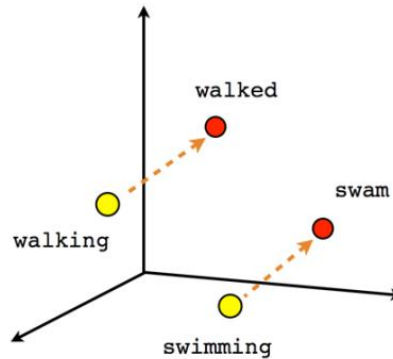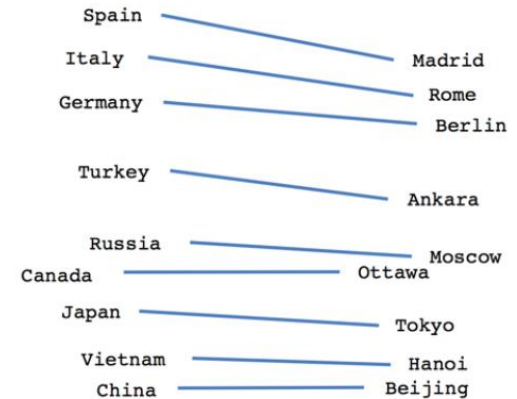Russia — Moscow

Canada — Ottawa

Japan — Tokyo

Vietnam — Hanoi

China — Beijing

Country-Capital

# Lecture 5: Text Representation II Distributed Representations

Pilsung Kang

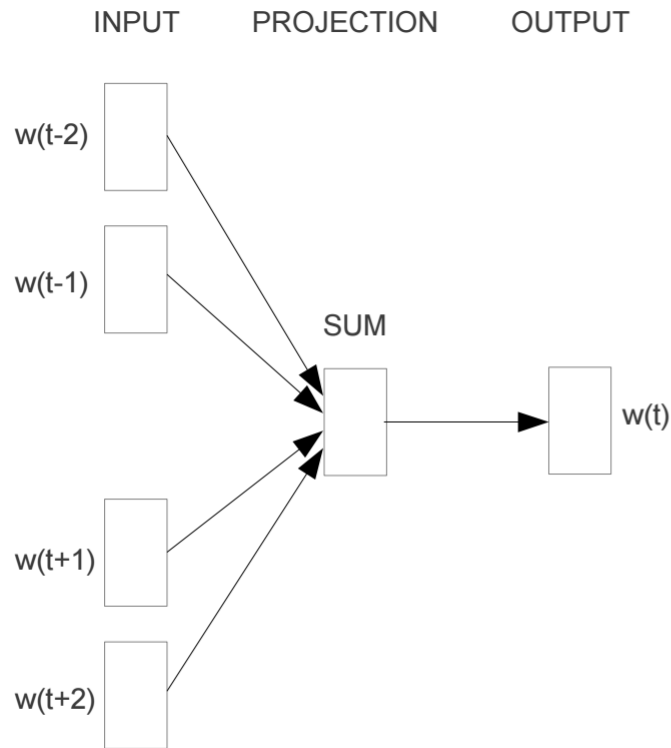School of Industrial Management Engineering

Korea University

# AGENDA

# Word2Vec

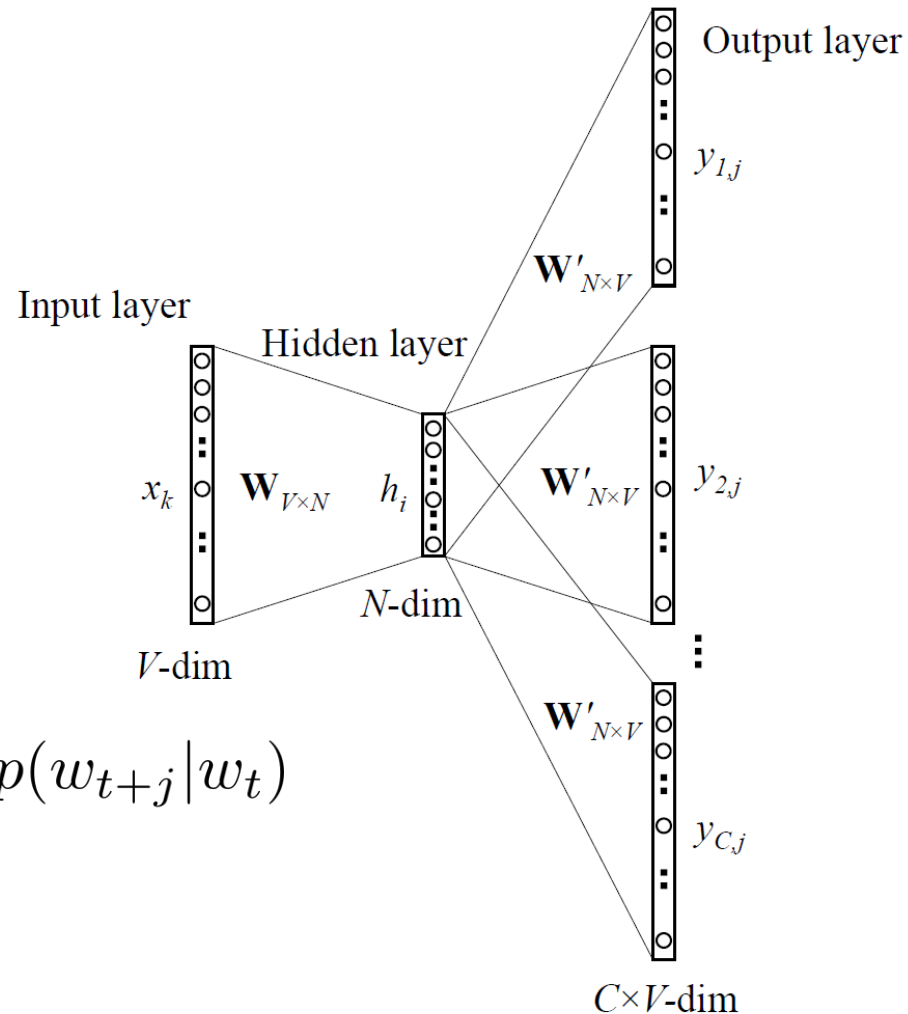- Two Architectures
  - ✓ Continuous bag-of-words (CBOW) vs. Skip-gram

# Word2Vec

- Learning representations: Skip-gram approach

  ✓ Predict surrounding words in a window of length m of every word

- Objective function

  ✓ Maximize the log probability of any context word given the current center word

$$J(\theta) = \frac{1}{T} \sum_{t=1}^{T} \sum_{-m \leq j \leq m, j \neq 0} \log p(w_{t+j}|w_t)$$
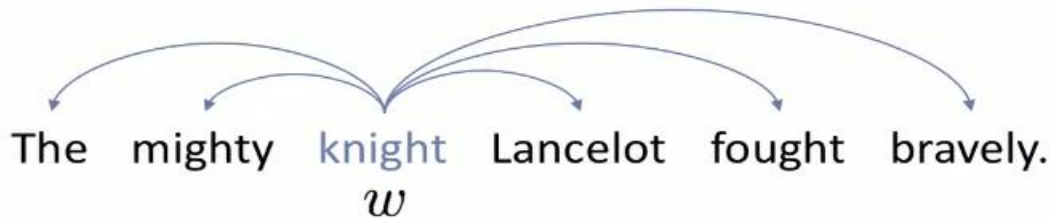
  ✓ where $\theta$ represents all variables we optimize

# Word2Vec

- Skip-gram model



knight ⟶ The
knight ⟶ mighty
knight ⟶ Lancelot
knight ⟶ fought
knight ⟶ bravely.

The  mighty  knight  Lancelot  fought  bravely.
$w$

- Model probability of a context word given a word
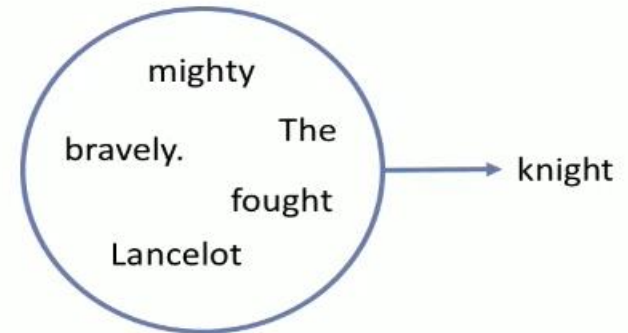
feature for word $w$: $x_w$
classifier for word $c$: $v_c$

$$p(c|w) = \frac{e^{x_w^\top v_c}}{\sum_{k=1}^{K} e^{x_w^\top v_k}}$$

- Word vectors $\quad x_w \in \mathbb{R}^d$

# Word2Vec

- CBOW model



- Model probability of a word given a context

feature for context $\mathcal{C}$: $h_{\mathcal{C}}$
classifier for word $w$: $v_w$

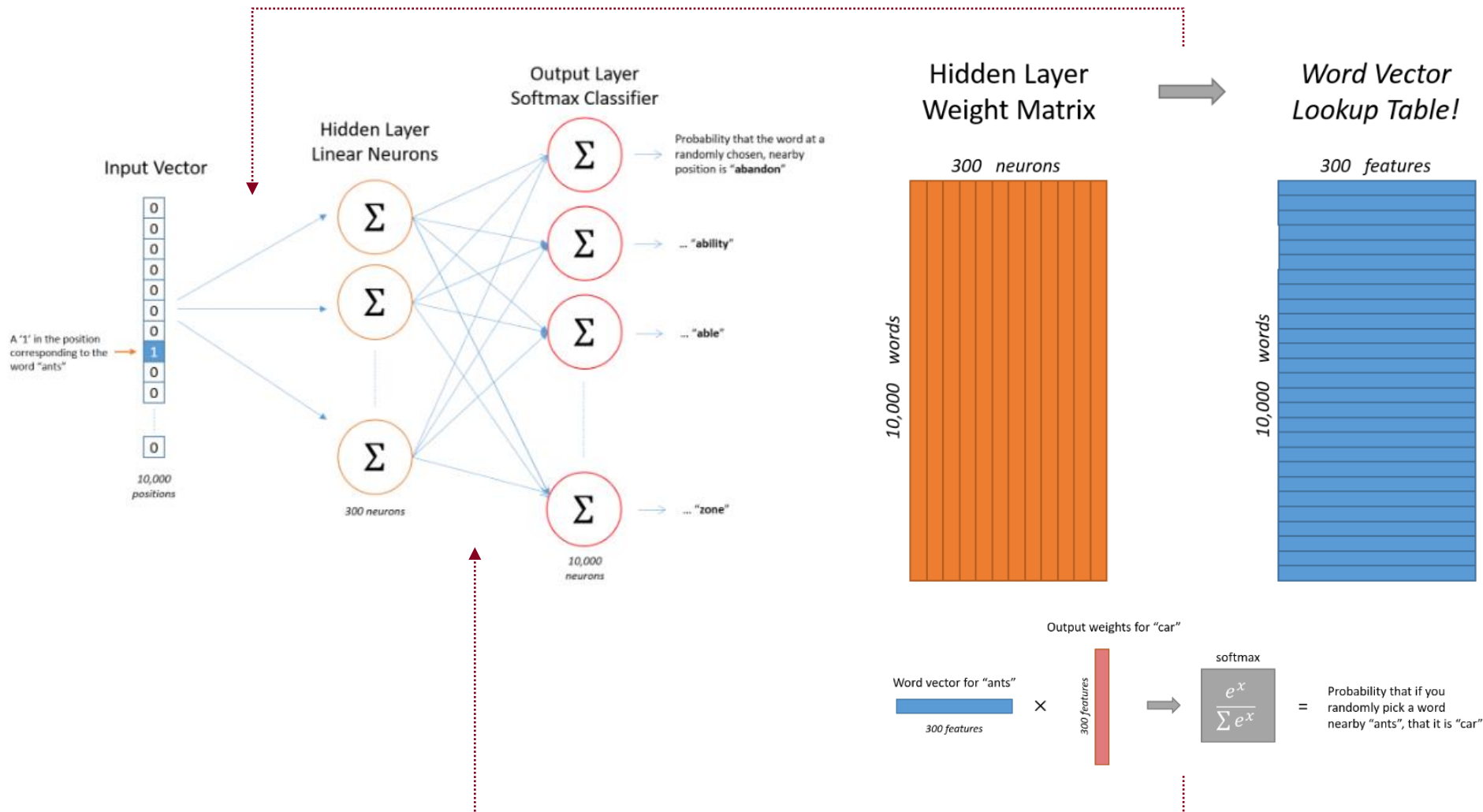$$p(w|\mathcal{C}) = \frac{e^{h_{\mathcal{C}}^{\top} v_w}}{\sum_{k=1}^{K} e^{h_{\mathcal{C}}^{\top} v_k}}$$

- Continuous Bag Of Words

$$h_{\mathcal{C}} = \sum_{c \in \mathcal{C}} x_c$$

# Word2Vec

- Another architecture explanation

# Word2Vec

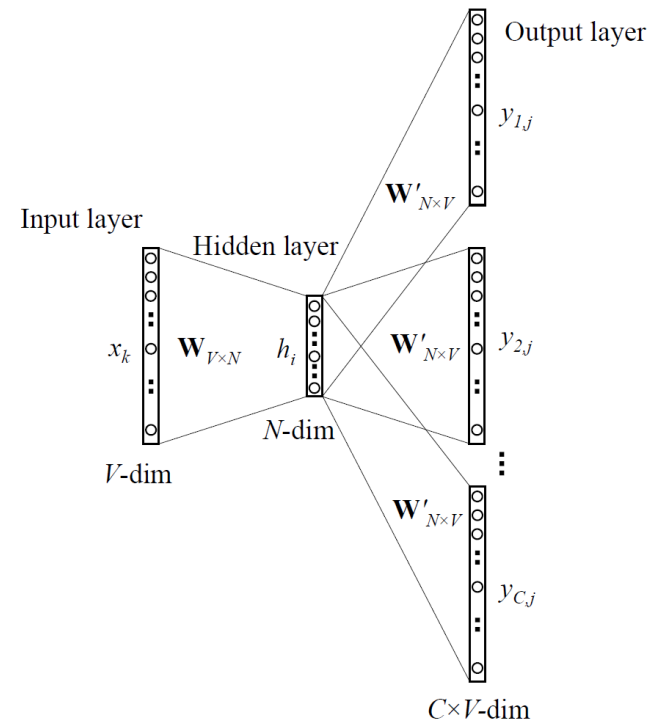- For simplicity, we use the following notation instead of $p(w_{t+j}|w_t)$

$$p(o|c) = \frac{exp(u_o^T v_c)}{\sum_{w=1}^{W} exp(u_w^T v_c)}$$

✓ where o is the outside (output) word id, c is the center word id, u and v are "outside" and "center" vectors of o and c

✓ Every word has two vectors!

- v is a row of matrix W

- u is a column of matrix W'

✓ Use W' = W$^T$ in practice for efficient computation

# Word2Vec

- Learning parameters with Gradient Ascent

$$J(\theta) = \frac{1}{T} \sum_{t=1}^{T} \sum_{-m \leq j \leq m, j \neq 0} \log p(w_{t+j}|w_t)$$

$$p(o|c) = \frac{exp(u_o^T v_c)}{\sum_{w=1}^{W} exp(u_w^T v_c)}$$

✓ Compute the gradient

$$\frac{\partial}{\partial v_c} \log p(o|c) = \frac{\partial}{\partial v_c} \log \frac{exp(u_o^T v_c)}{\sum_{w=1}^{W} exp(u_w^T v_c)}$$

$$= \frac{\partial}{\partial v_c} u_o^T v_c - \frac{\partial}{\partial v_c} \log \sum_{w=1}^{W} exp(u_w^T v_c)$$

A          B

# Word2Vec

- Learning parameters with Gradient Ascent

  ✓ For chunk A

  $$\frac{\partial}{\partial v_c} u_o^T v_c = u_o$$

  ✓ For chunk B

  $$-\frac{\partial}{\partial v_c} \log \sum_{w=1}^{W} exp(u_w^T v_c)$$

  $$= -\frac{1}{\sum_{w=1}^{W} exp(u_w^T v_c)} \cdot \left( \sum_{w=1}^{W} exp(u_w^T v_c) \cdot u_w \right)$$

  $$= -\sum_{w=1}^{W} \frac{exp(u_w^T v_c)}{\sum_{w=1}^{W} exp(u_w^T v_c)} \cdot u_w = -\sum_{w=1}^{W} P(w|c) \cdot u_w$$

# Word2Vec

- Learning parameters with Gradient Ascent

$$\frac{\partial}{\partial v_c}\log p(o|c) = u_o - \sum_{w=1}^{W} P(w|c) \cdot u_w$$

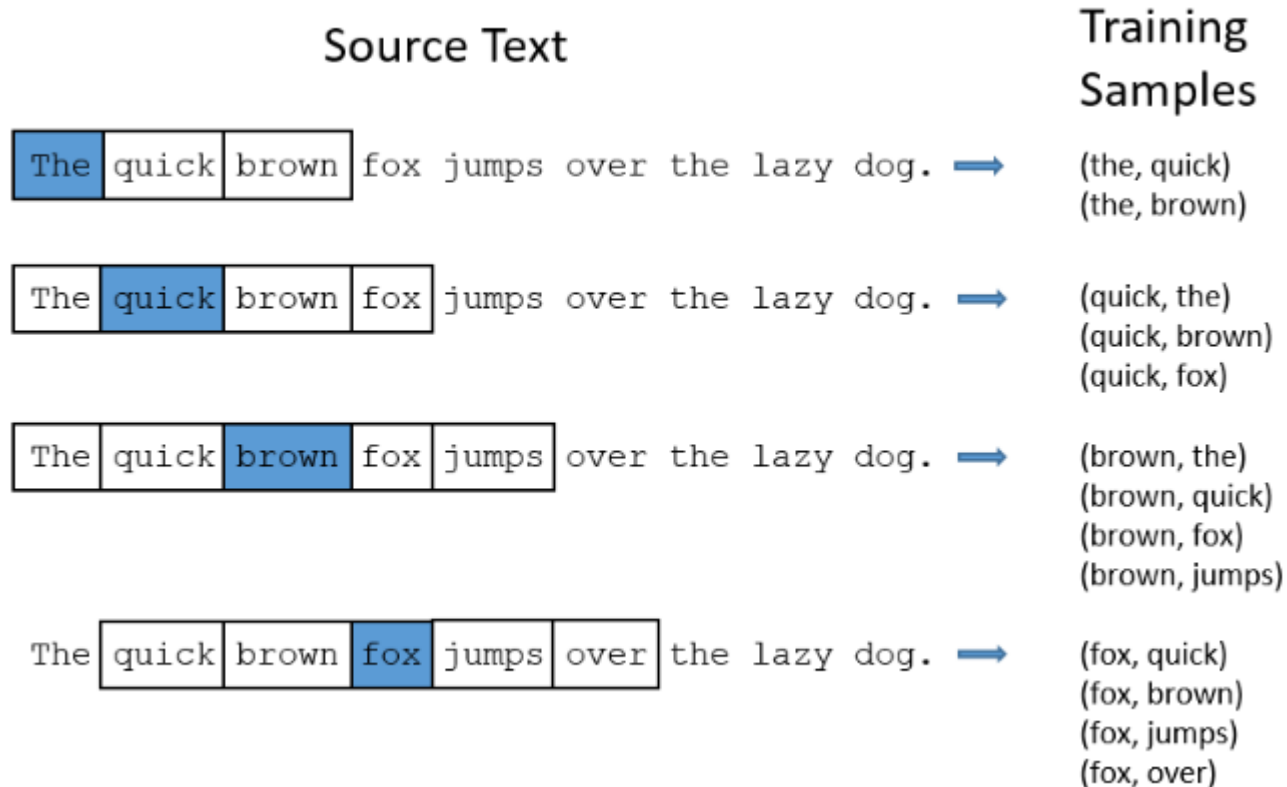- Update the weight vector

$$v_c(t+1) = v_c(t) + \alpha\left(u_o - \sum_{w=1}^{W} P(w|c) \cdot u_w\right)$$

# Word2Vec

- Learning strategy

  ✓ Do not use all nearby words, but one per each training

### Source Text

### Training Samples

The **quick** brown fox jumps over the lazy dog. ➡ (the, quick)
(the, brown)

The **quick** brown fox jumps over the lazy dog. ➡ (quick, the)
(quick, brown)
(quick, fox)

The quick **brown** fox jumps over the lazy dog. ➡ (brown, the)
(brown, quick)
(brown, fox)
(brown, jumps)

The quick brown **fox** jumps over the lazy dog. ➡ (fox, quick)
(fox, brown)
(fox, jumps)
(fox, over)

# Word2Vec

- The number of weights to be trained: 2 x V x N (Huge network!)

  - ✓ Word pairs and phrases

    - ▪ Treating common word pairs or phrases as single "word"

  - ✓ Subsampling frequent words

    - ▪ To decrease the number of training examples

    - ▪ The probability of word $w_i$ being removed

$$P(w_i) = 1 - \sqrt{\frac{t}{f(w_i)}}$$

⟵ Threshold ($10^{-5}$)

⟵ Term frequency in the corpus

$$if\ f(w_i) = 10^{-4},\ P(w_i) = 1 - \sqrt{\frac{1}{10}} = 0.6838$$

$$if\ f(w_i) = 10^{-2},\ P(w_i) = 1 - \sqrt{\frac{1}{1000}} = 0.9684$$

# Word2Vec

- The number of weights to be trained: 2 x V x N (Huge network!)

  ✓ Negative sampling

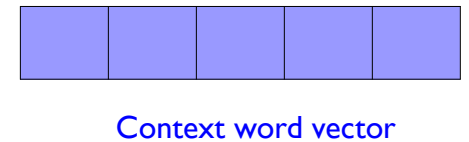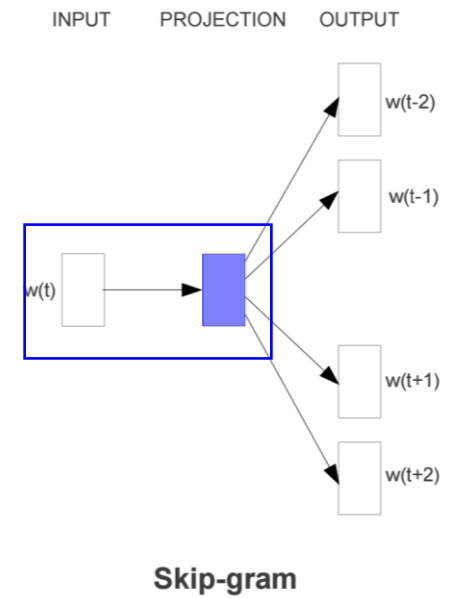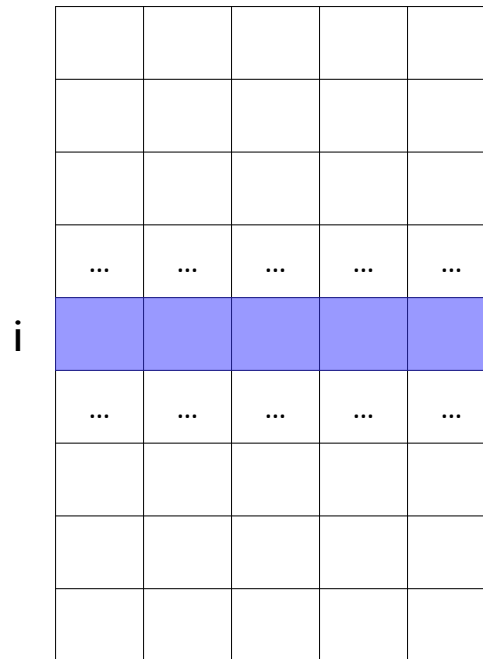    ▪ Instead of updating the weights associated with all output words, update the weight of a few (5-20) words
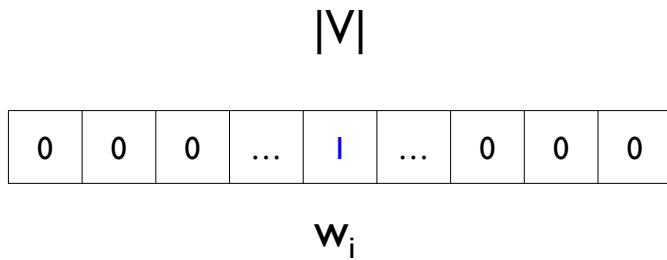
$$J(\theta) = \frac{1}{T}\sum_{t=1}^{T}\sum_{-m\leq j\leq m, j\neq 0} \log p(w_{t+j}|w_t) = \frac{1}{T}\sum_{t=1}^{T} J_t(\theta)$$

$$J_t(\theta) = \log\sigma(u_o^T v_c) + \sum_{i=1}^{k} E_{i\sim P(w)}\left[\log\sigma(-u_i^T v_c)\right]$$

$$P(w_i) = \frac{f(w_i)^{3/4}}{\sum_{j=0}^{n}\left(f(w_i)^{3/4}\right)}$$

# Word2Vec

- Negative Sampling Example

$|V|$

| 0 | 0 | 0 | … | 1 | … | 0 | 0 | 0 |

$w_i$



INPUT    PROJECTION    OUTPUT
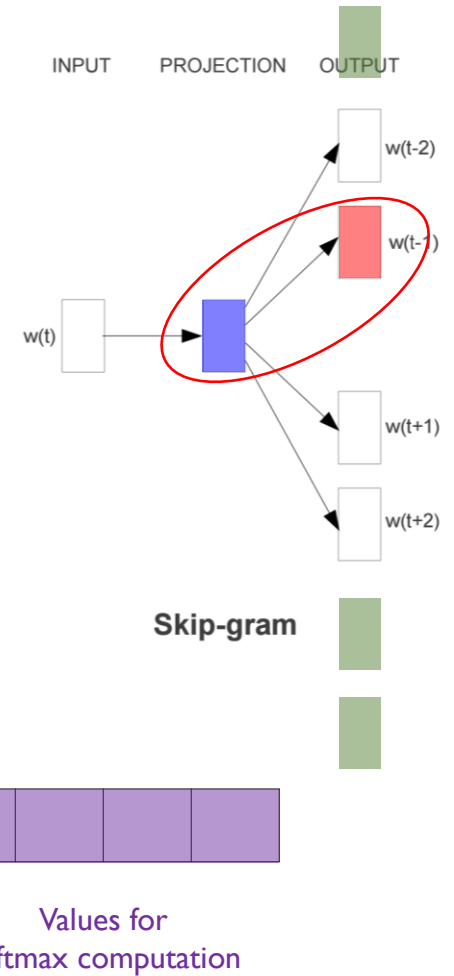
w(t-2)
w(t-1)
w(t)
w(t+1)
w(t+2)

**Skip-gram**

i

Context word vector
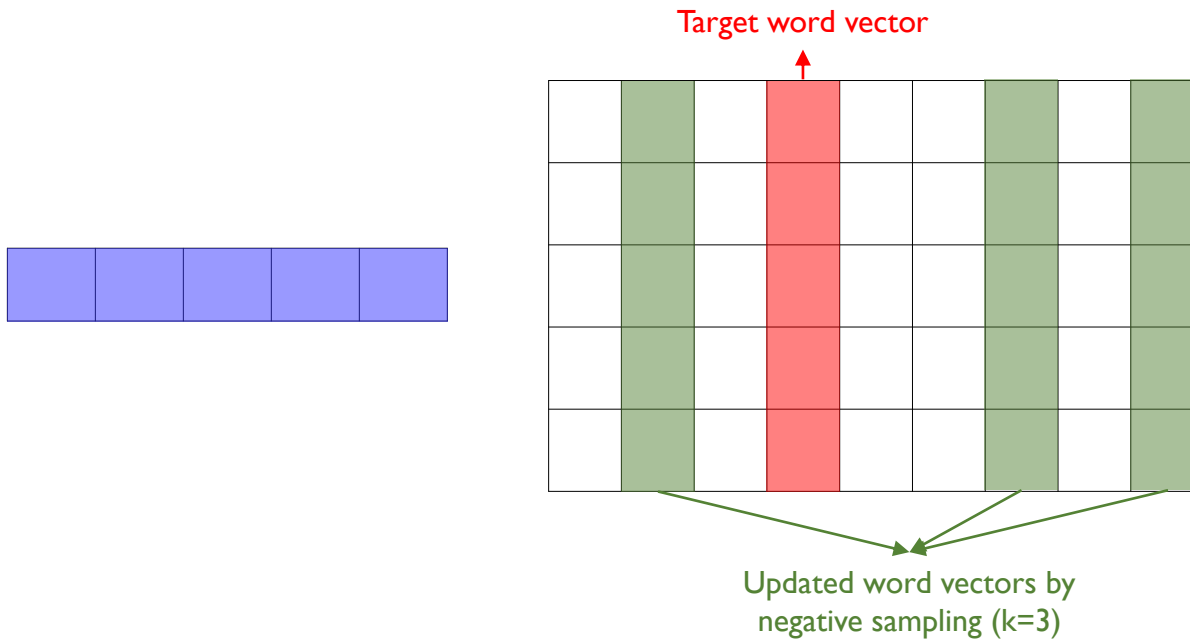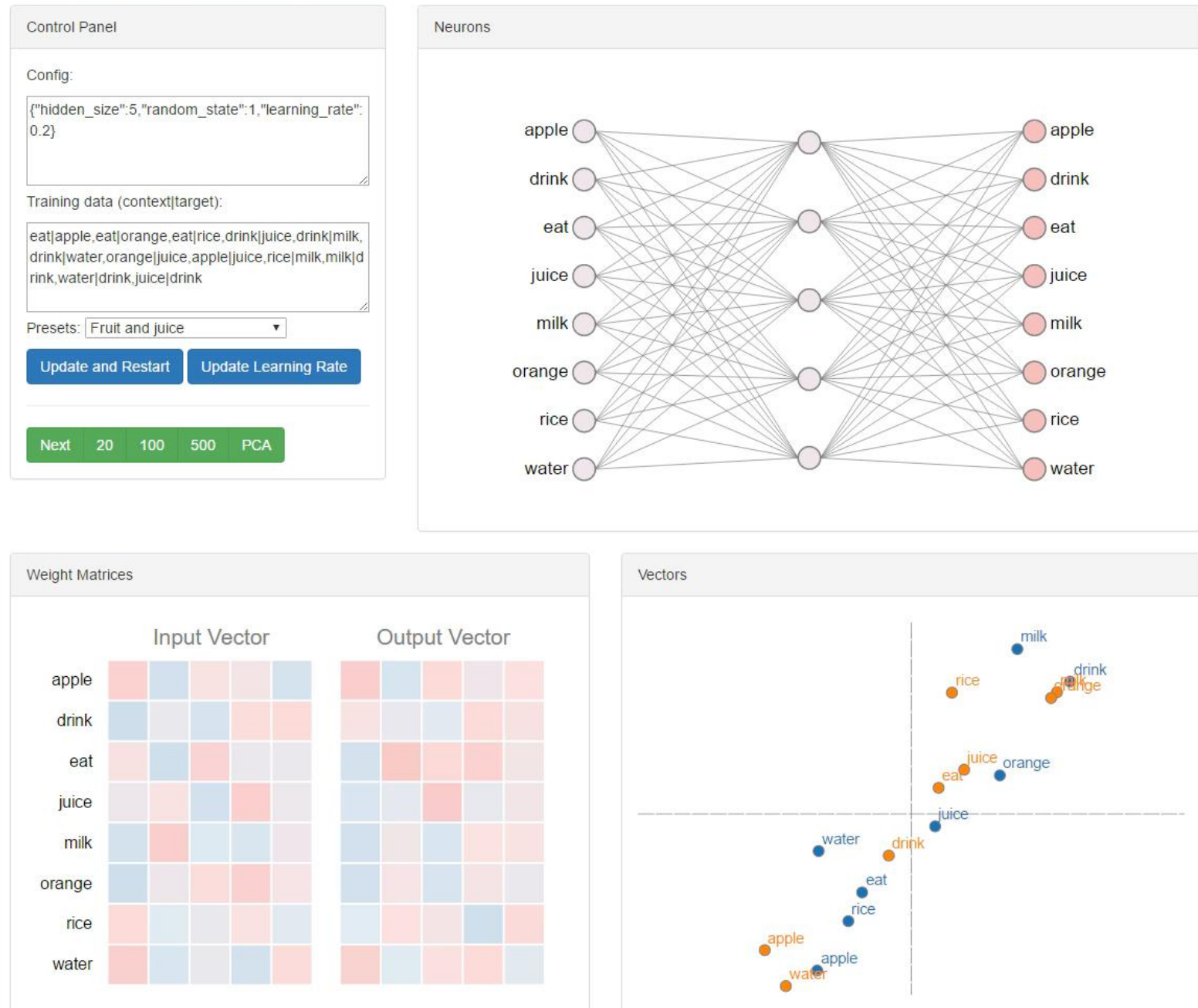
# Word2Vec

- Negative Sampling Example

  ✓ No. of 0 values for softmax computation is reduced from

  |V| to (k+1)

INPUT    PROJECTION    OUTPUT

w(t)

w(t-2)

w(t-1)

w(t+1)

w(t+2)

**Skip-gram**

Target word vector

Updated word vectors by
negative sampling (k=3)
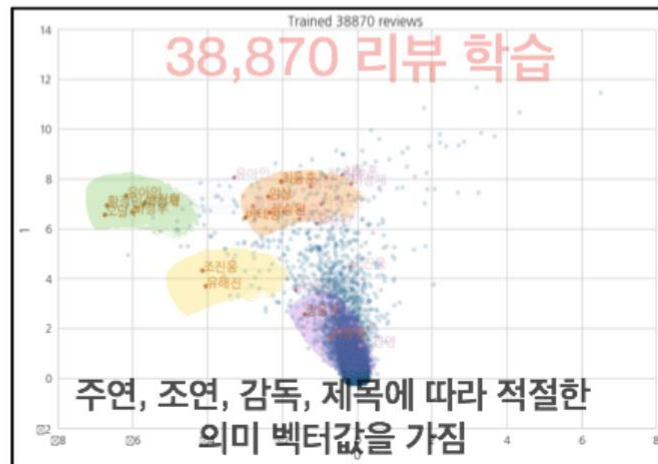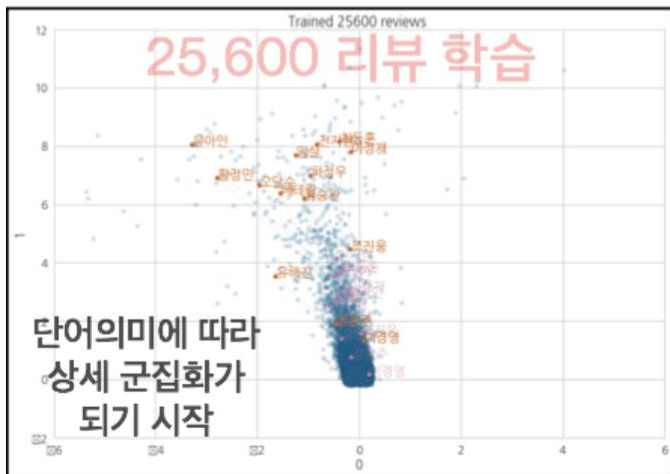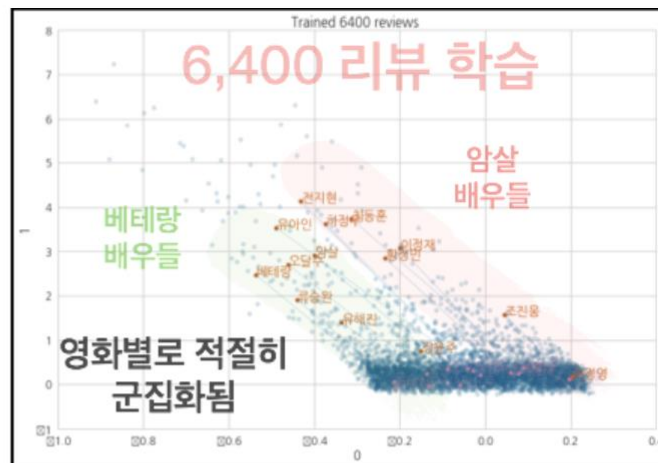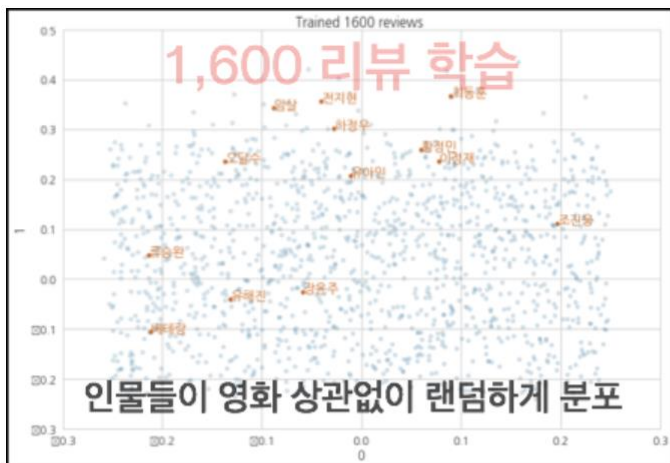
Values for
softmax computation

# wevi: word embedding visual inspector

Everything you need to know about this tool - Source code

## Control Panel

Config:

```
{"hidden_size":5,"random_state":1,"learning_rate":
0.2}
```

Training data (context|target):

```
eat|apple,eat|orange,eat|rice,drink|juice,drink|milk,
drink|water,orange|juice,apple|juice,rice|milk,milk|d
rink,water|drink,juice|drink
```

Presets: Fruit and juice

[ Update and Restart ] [ Update Learning Rate ]

[ Next ] [ 20 ] [ 100 ] [ 500 ] [ PCA ]

## Neurons



apple, drink, eat, juice, milk, orange, rice, water

## Weight Matrices

### Input Vector

### Output Vector

apple, drink, eat, juice, milk, orange, rice, water

## Vectors



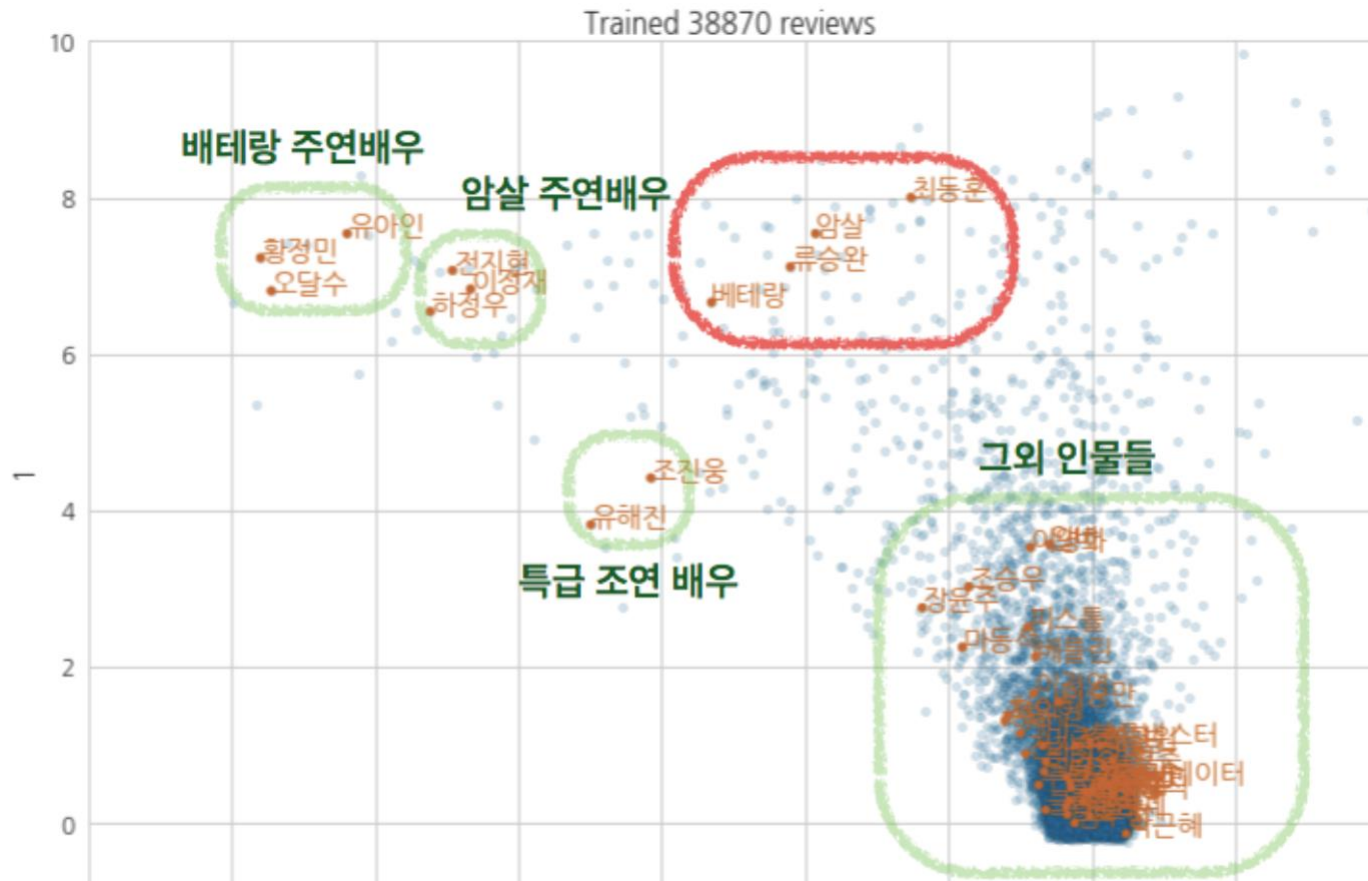milk, rice, drink, milk, orange, juice, orange, eat, juice, water, drink, eat, rice, apple, apple, water

# Word2Vec: Example

- Trained with movie reviews

# Word2Vec: Example

- Trained with movie reviews



$$Vec[류승완] - Vec[베테랑] \doteqdot Vec[최동훈] - Vec[암살]$$