

Paper Seminar

# **Analogies Explained: Towards Understanding Word Embeddings**

**Allen and Hospedales, 2019, ICML**

**Myeongsup Kim**

Integrated M.S./Ph.D. Student  
Data Science & Business Analytics Lab.  
School of Industrial Management Engineering  
Korea University

Myeongsup\_kim@korea.ac.kr

# Introduction

- Analogies in Word Embeddings

## Introduction

-What This Seminar Does Not Cover

### <What This Seminar Does Not Cover>

- **Details of Word2Vec**

[Mikolov et al., Efficient Estimation of Word Representations in Vector Space, ICLR Workshop, 2013](#)

[Mikolov et al., Distributed Representations of Words and Phrases and their Compositionality, NIPS, 2013](#)

## Introduction

-Analogies in Word Embeddings

### <Analogies in Natural Language>

“man is to boy as woman is to girl”

## Introduction

-Analogies in Word Embeddings

### <Analogies in Natural Language>

“man is to boy as woman is to girl”

man : boy = woman : girl  
(In Mathematical Expression)

## Introduction

-Analogies in Word Embeddings

### <Analogies in Natural Language>

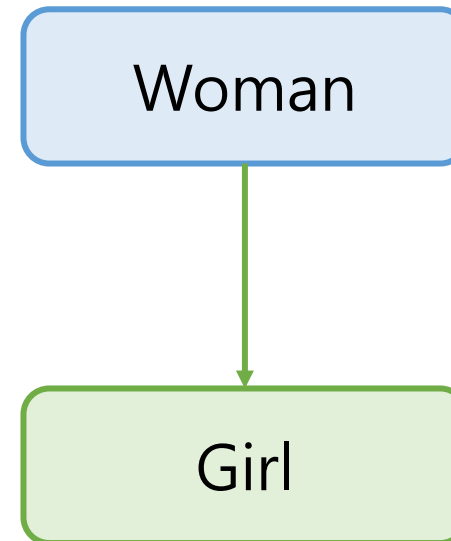
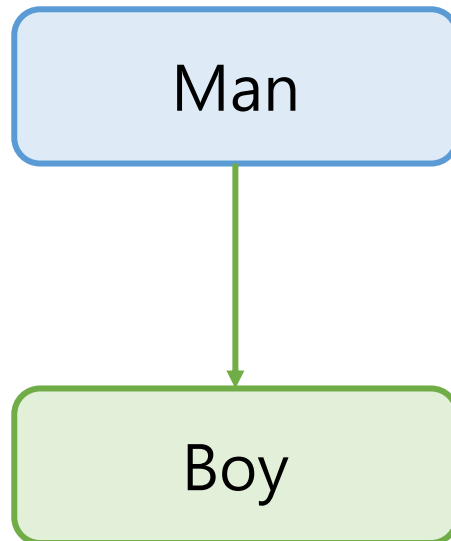
Man

Woman

## Introduction

-Analogies in Word Embeddings

### <Analogies in Natural Language>



## Introduction

-Analogies in Word Embeddings

### <Analogies in Natural Language>

"man is to king as woman is to ...?"

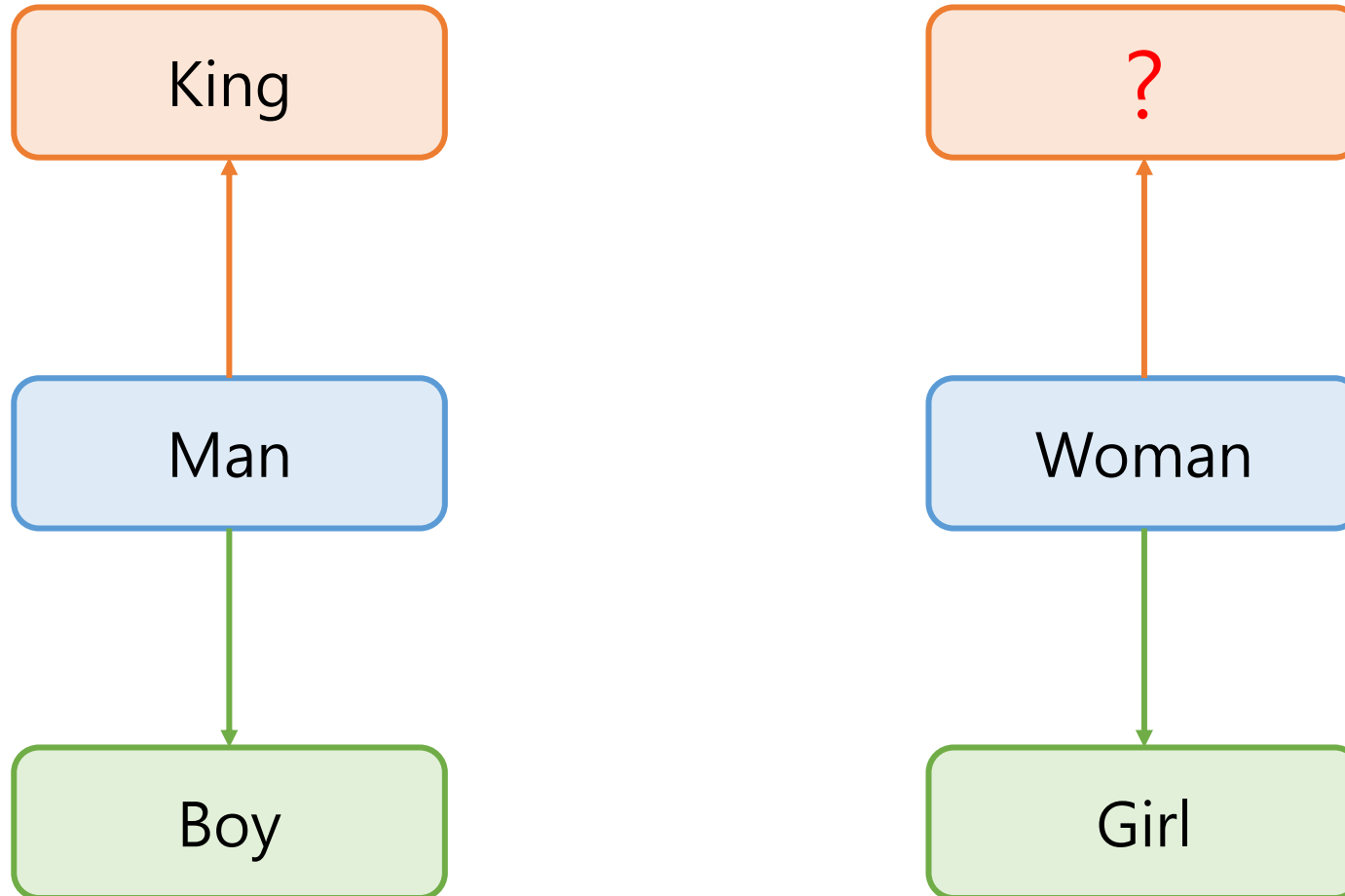
man : king = woman : ?  
(In Mathematical Expression)



## Introduction

-Analogies in Word Embeddings

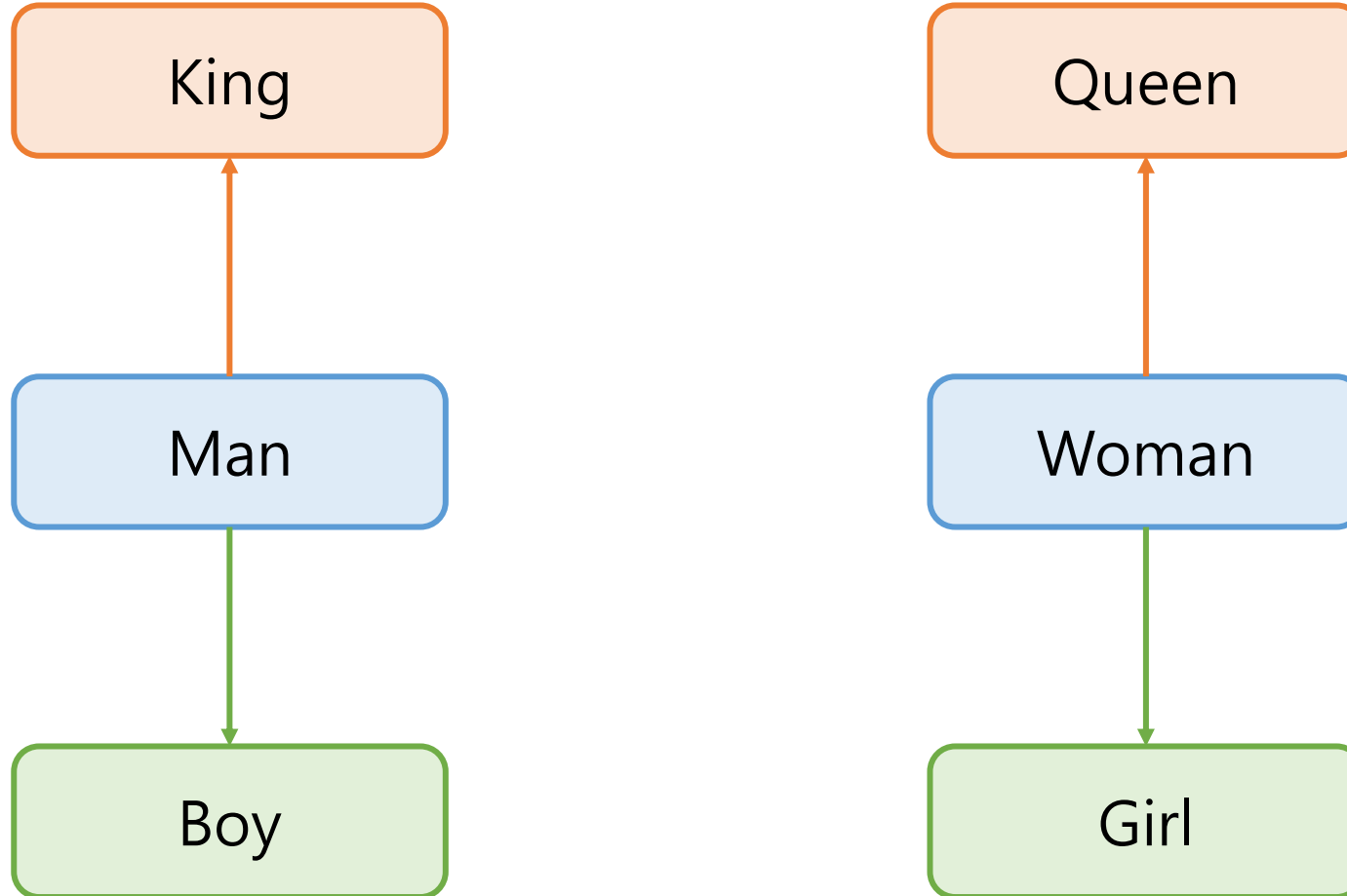
### <Analogies in Natural Language>



## Introduction

-Analogies in Word Embeddings

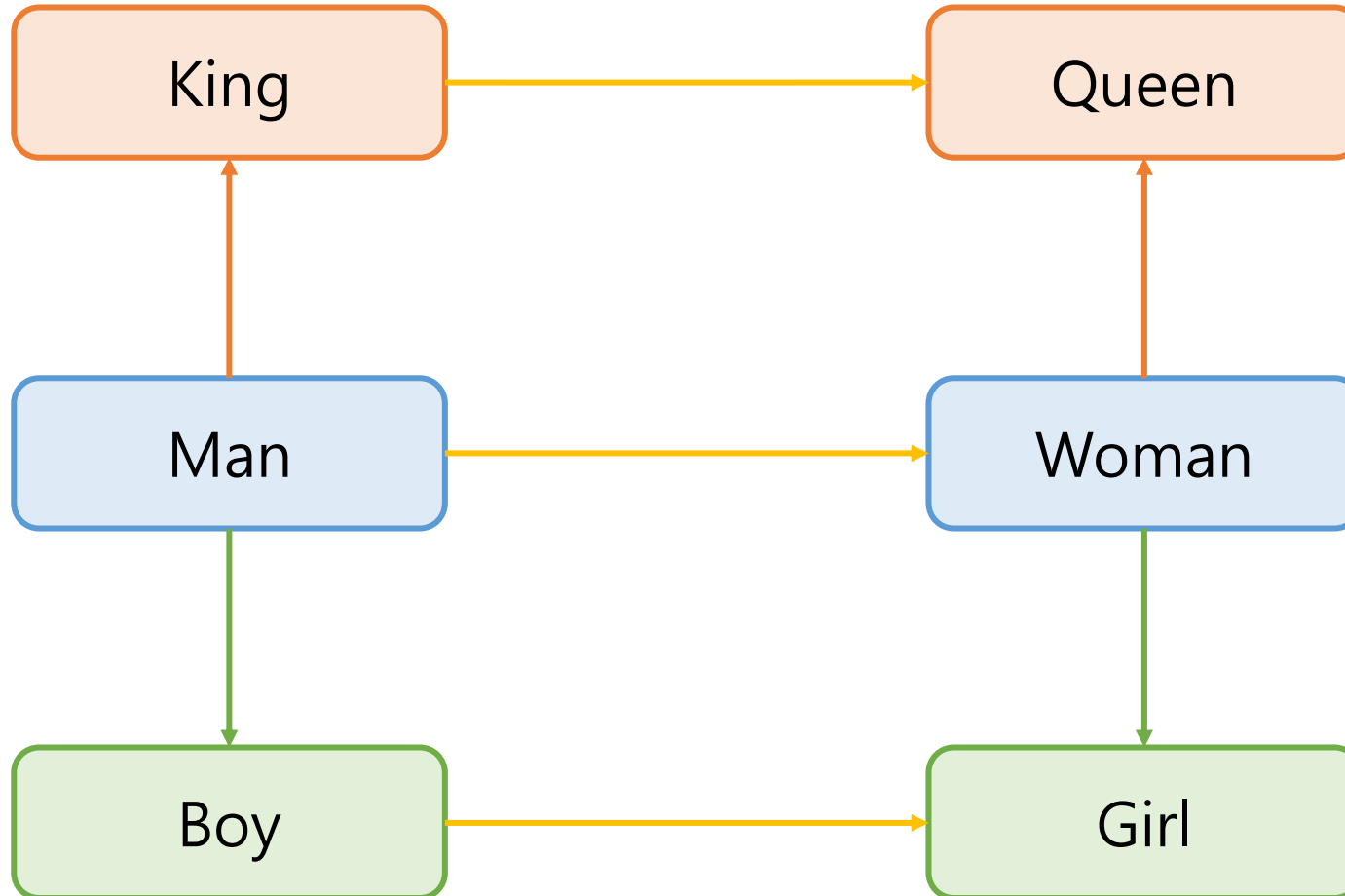
### <Analogies in Natural Language>



# Introduction

-Analogies in Word Embeddings

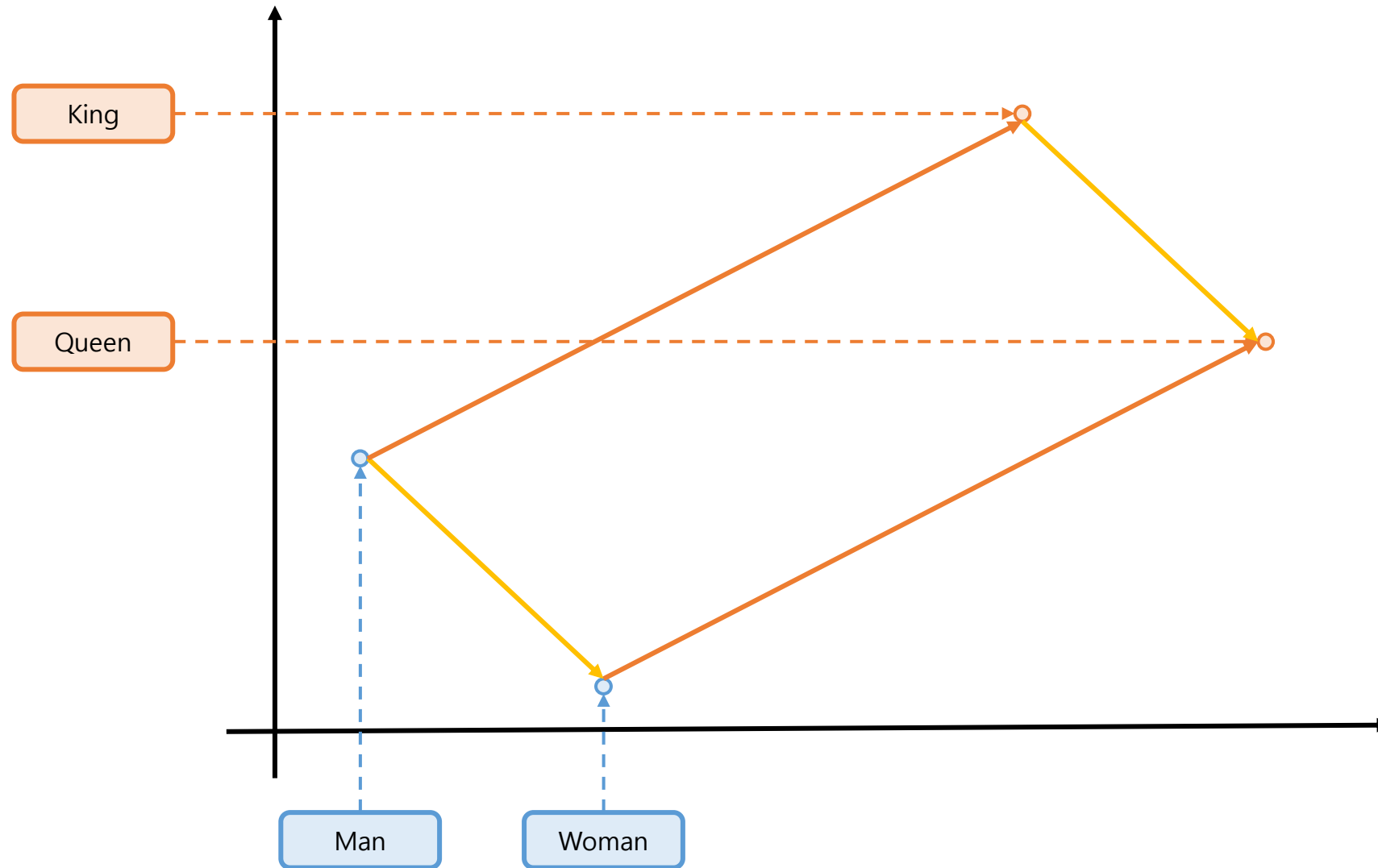
## <Analogies in Natural Language>



# Introduction

-Analogies in Word Embeddings

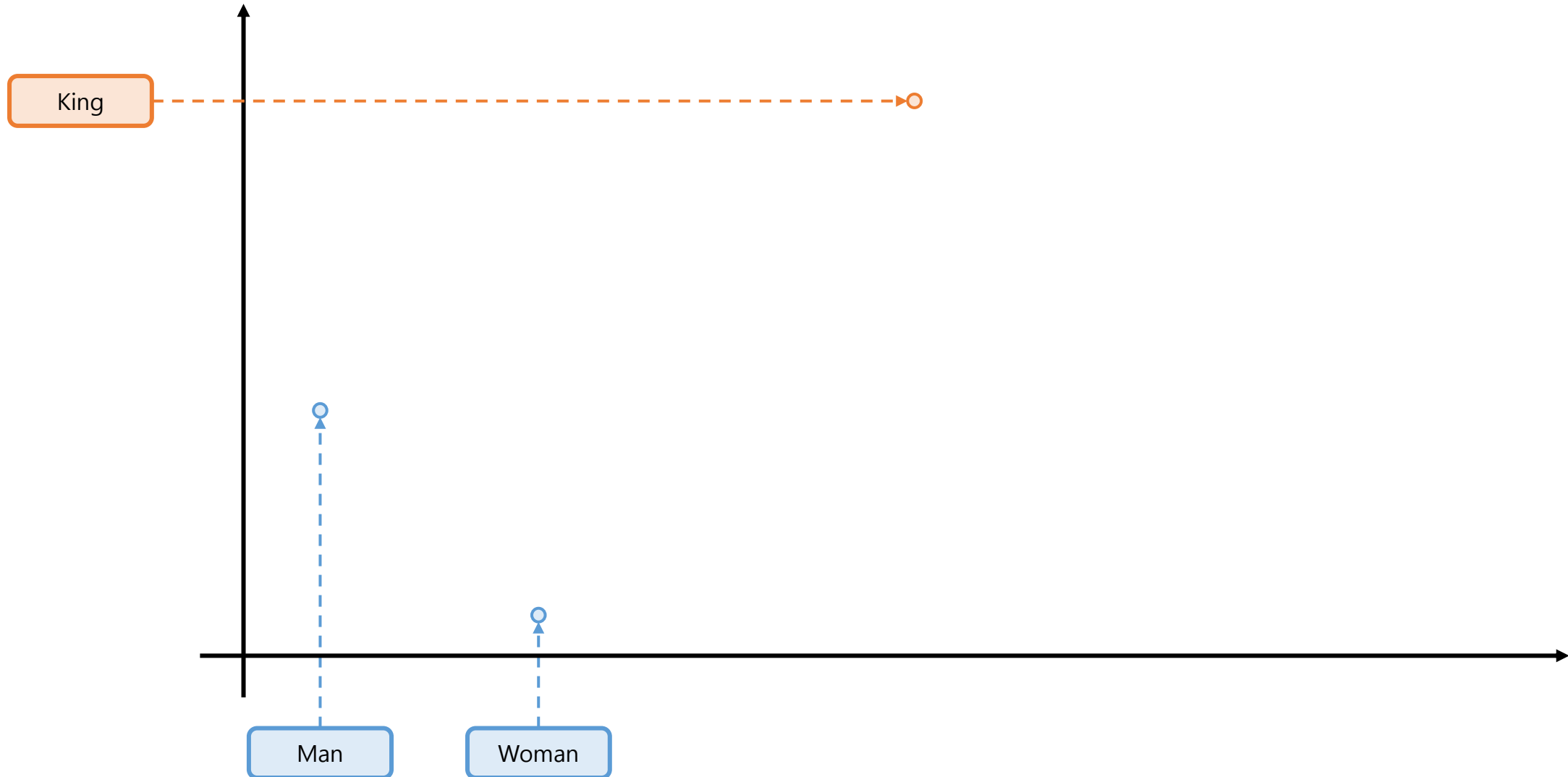
## <Analogies in Vector Representation>



# Introduction

-Analogies in Word Embeddings

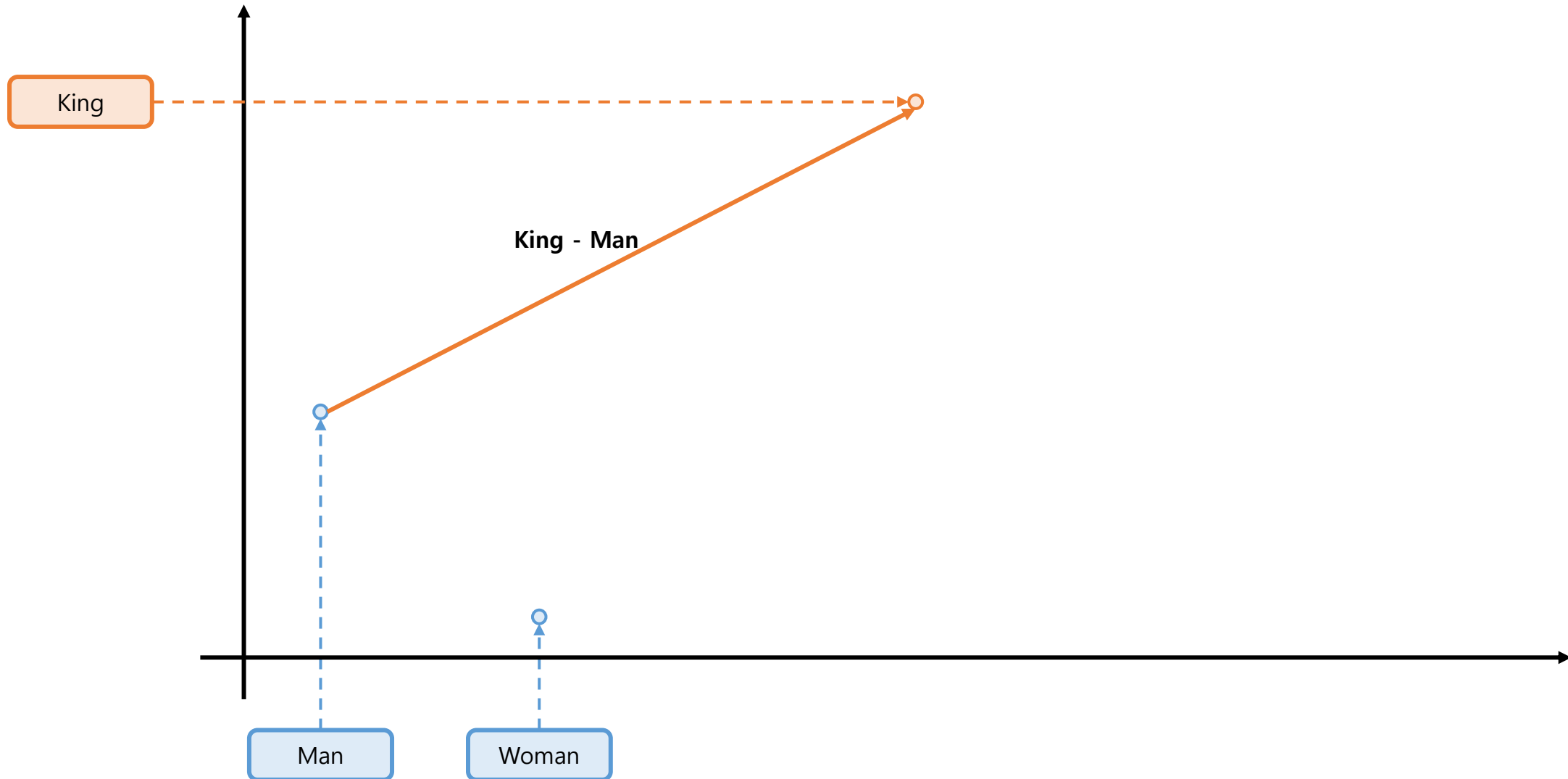
## <Analogies in Vector Representation>



# Introduction

-Analogies in Word Embeddings

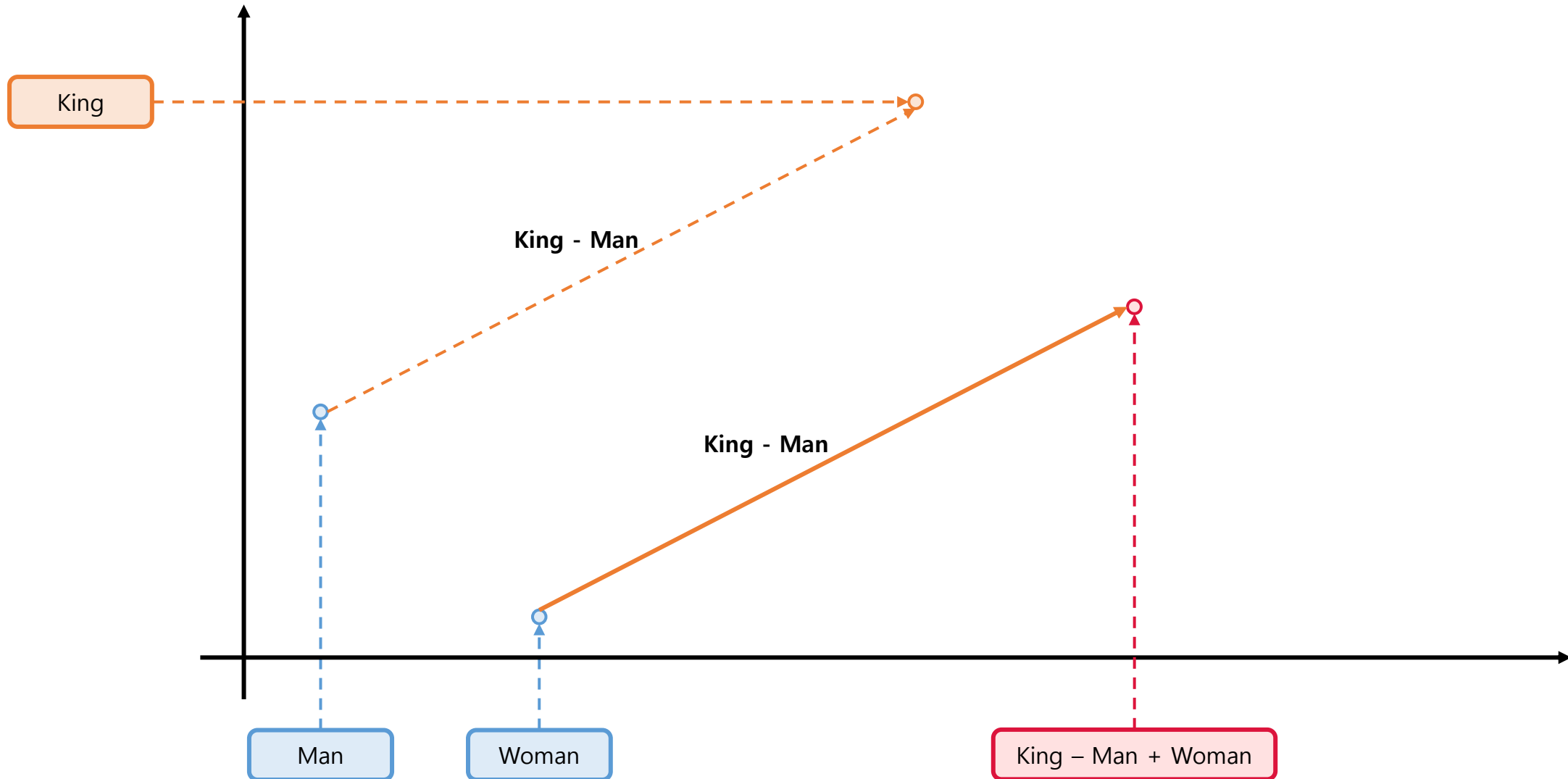
## <Analogies in Vector Representation>



# Introduction

-Analogies in Word Embeddings

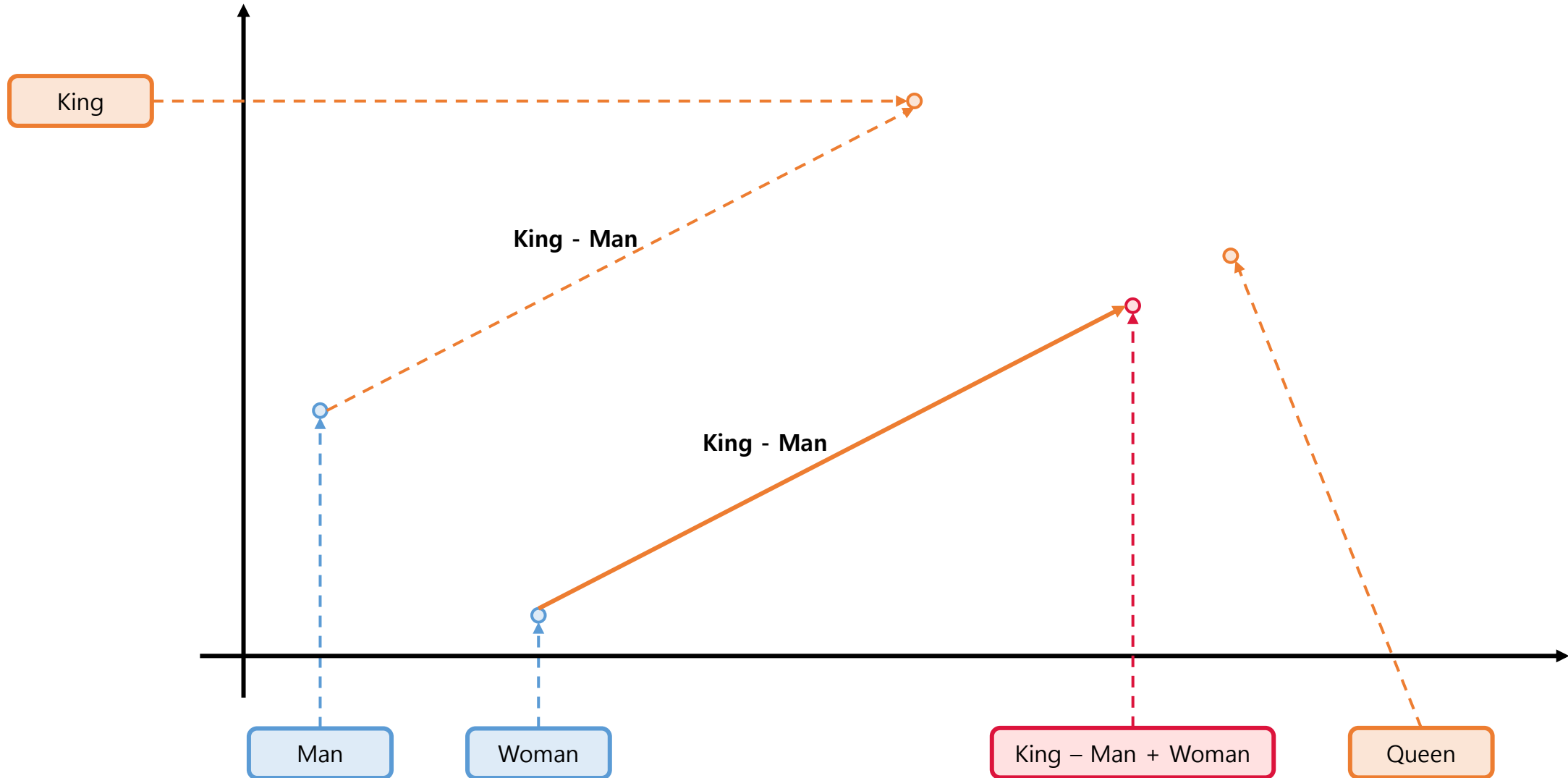
## <Analogies in Vector Representation>



# Introduction

-Analogies in Word Embeddings

## <Analogies in Vector Representation>

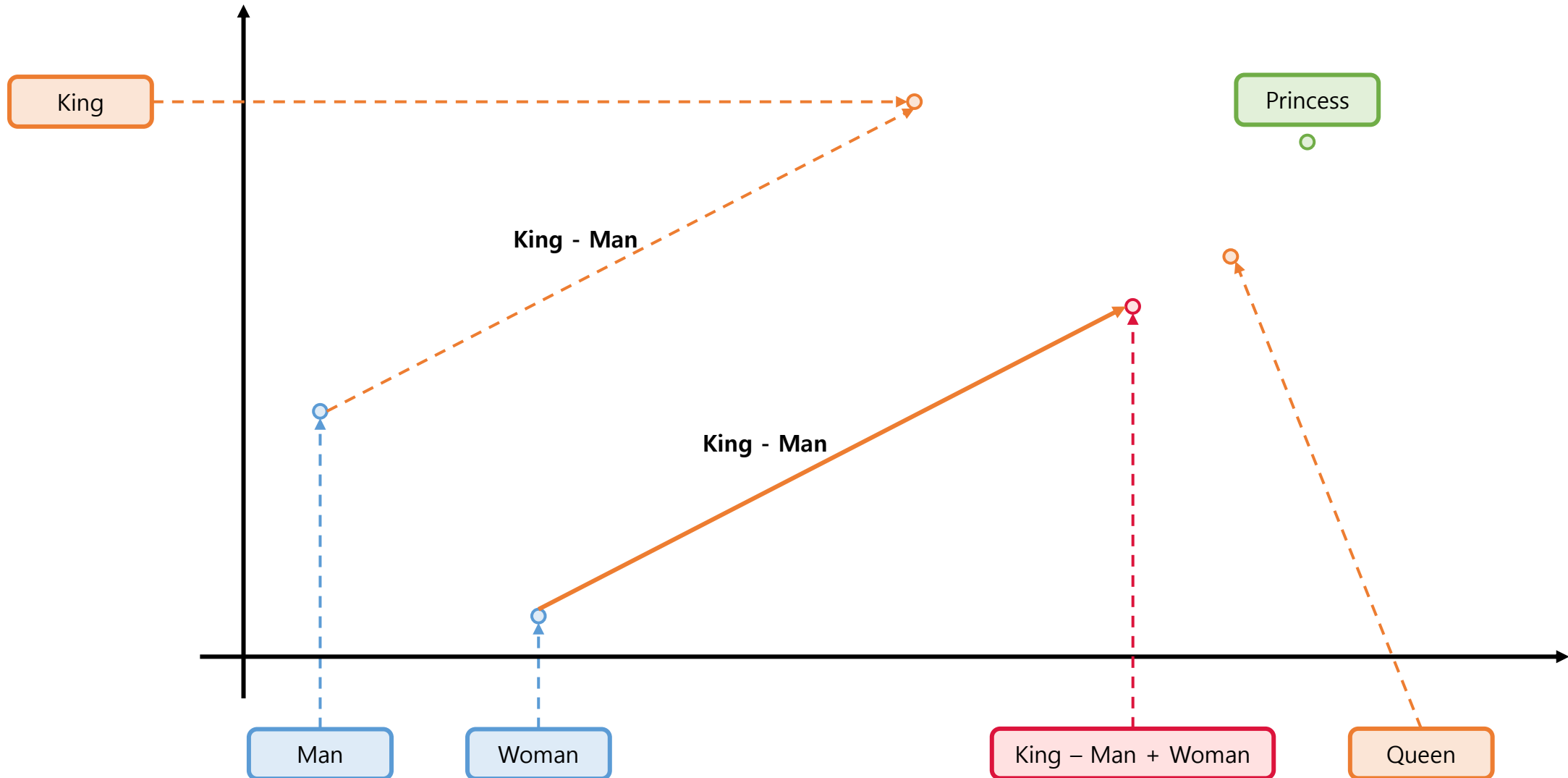




# Introduction

-Analogies in Word Embeddings

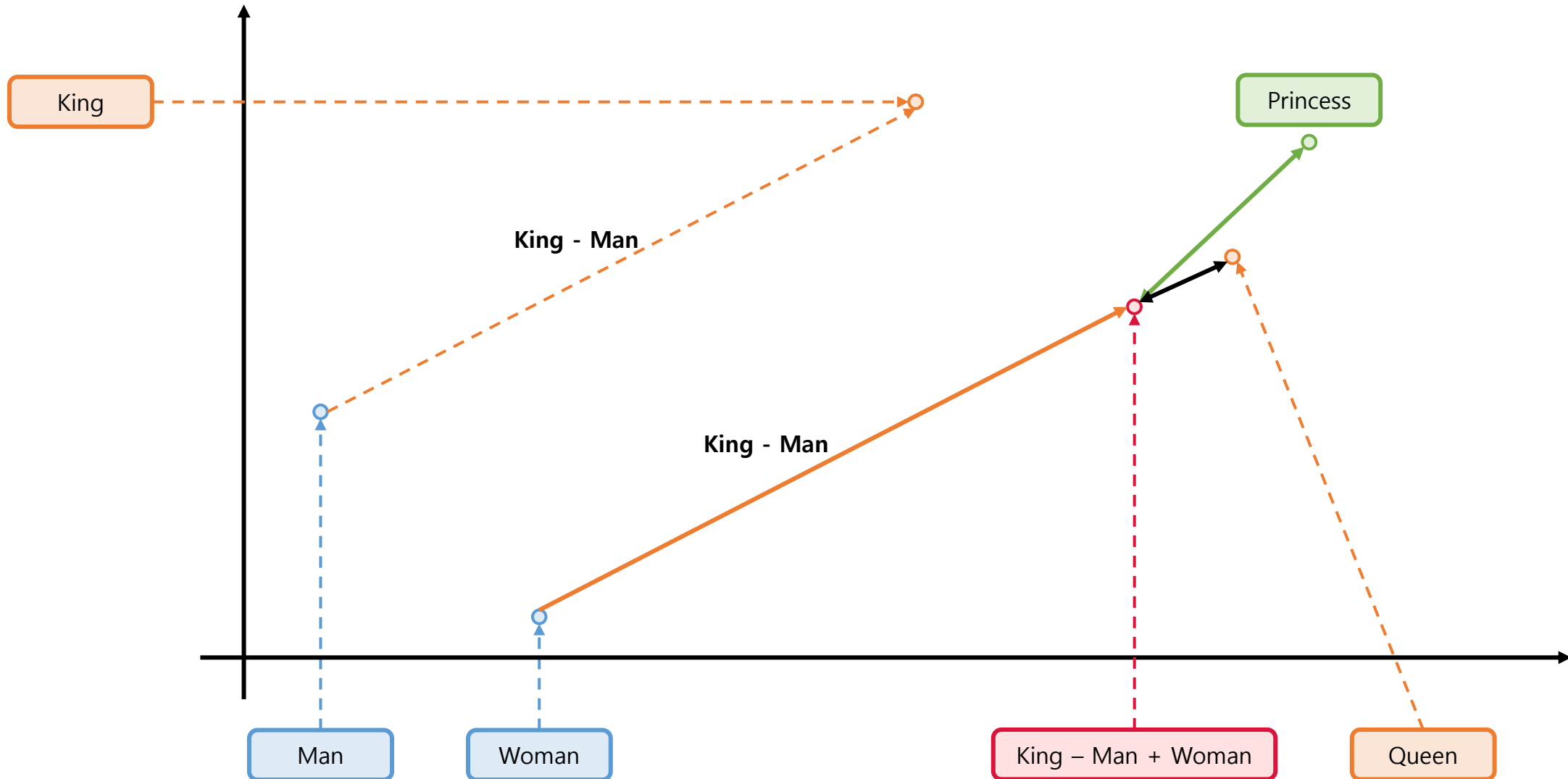
## <Analogies in Vector Representation>



# Introduction

-Analogies in Word Embeddings

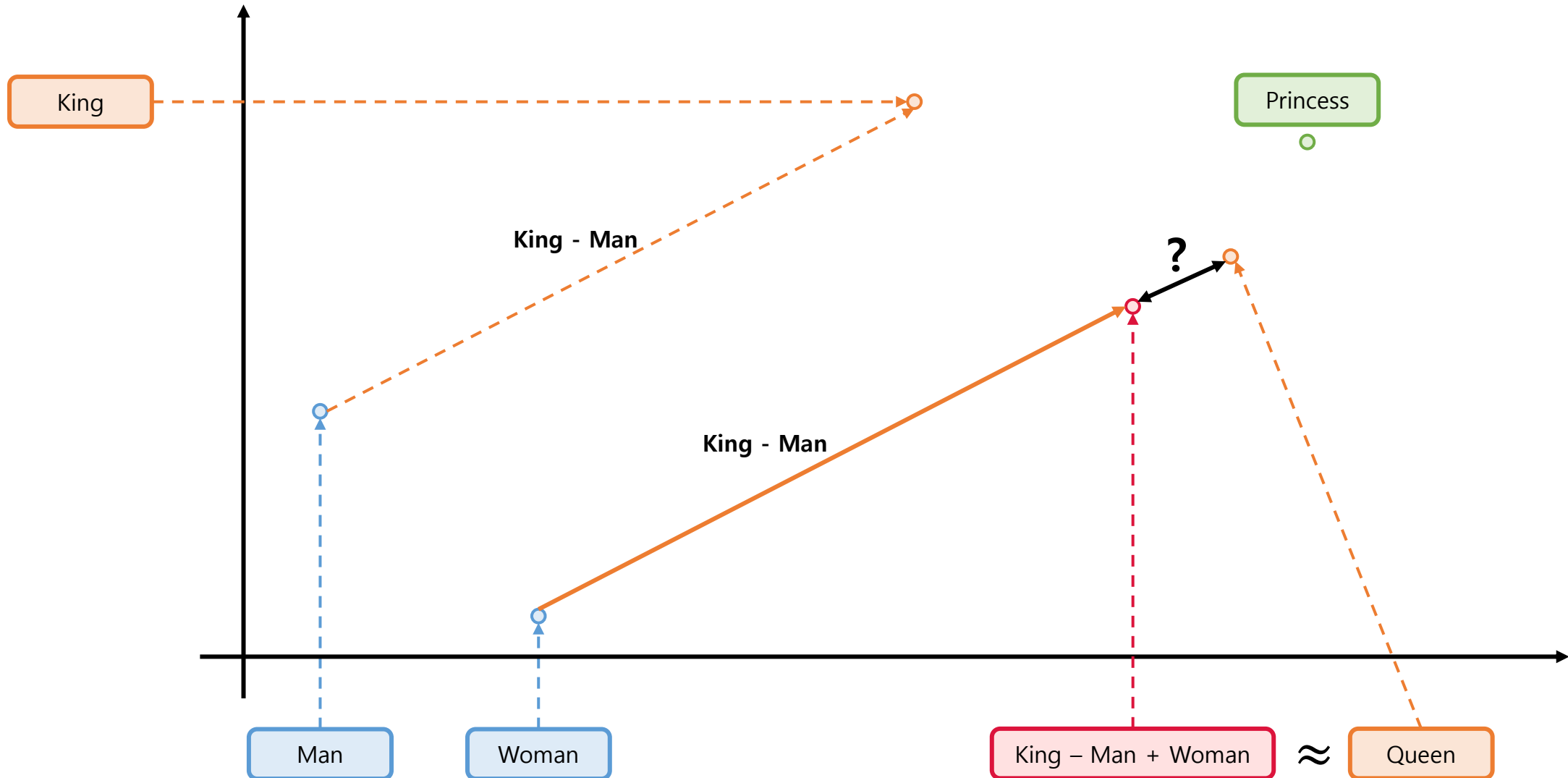
## <Analogies in Vector Representation>



# Introduction

-Analogies in Word Embeddings

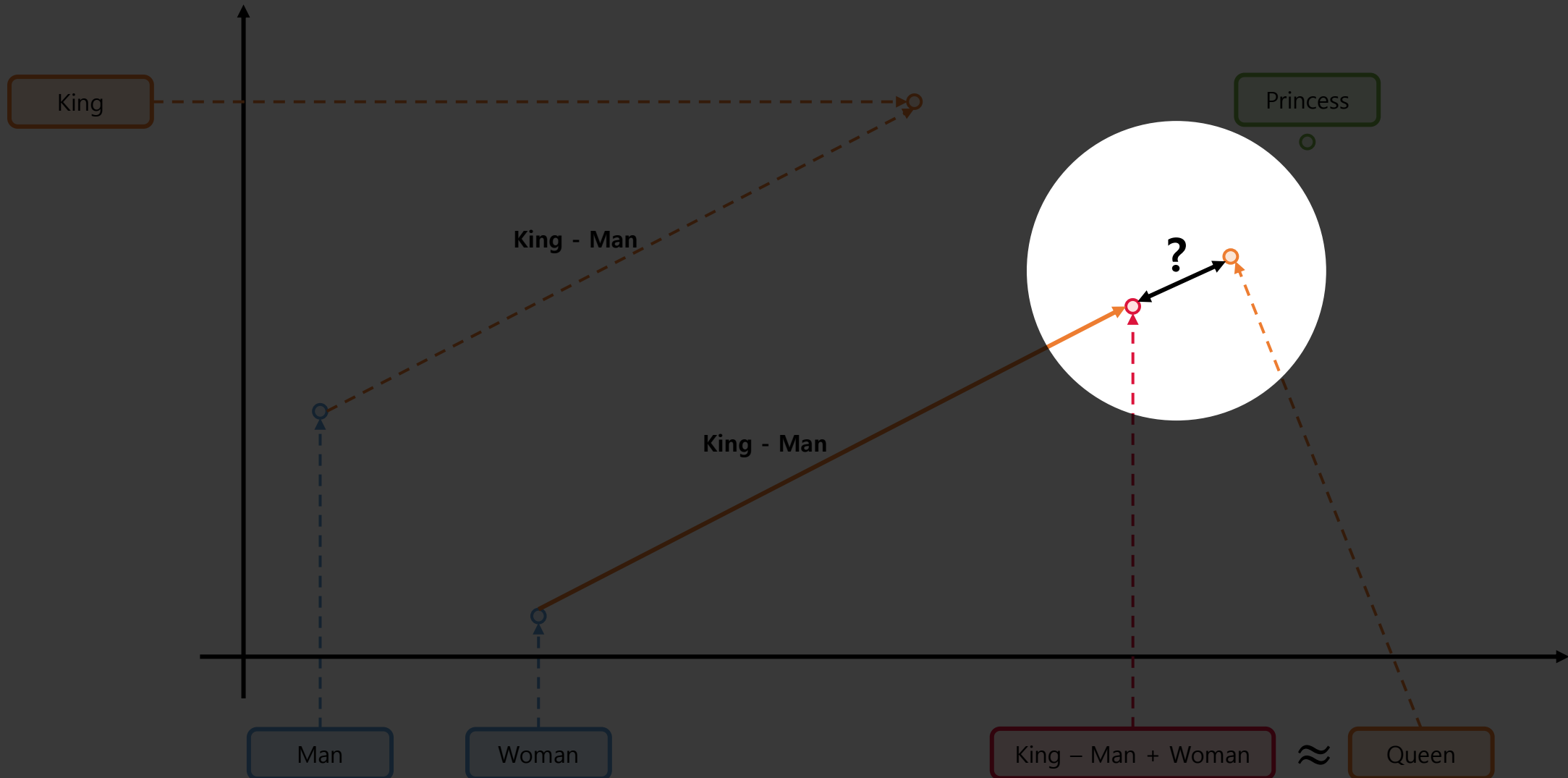
## <Analogies in Vector Representation>



# Introduction

-Analogies in Word Embeddings

## <Analogies in Vector Representation>



## -Analogies in Word Embeddings

King

Princess

King - Man

Man

Woman

King - Man + Woman

Queen

Why Does Analogy of Words Hold in Embedding Space, Even Though We haven't Trained in That Way?

Can We Interpret Analogy and Explain Linear Relationship Between Word Embeddings in a Mathematical Way?

# Can We Interpret Analogy and Explain Linear Relationship Between Word Embeddings in a Mathematical Way?

# Analogy Explained: Towards Understanding Word Embeddings

*Allen and Hospedales, 2019, ICML*

*"ICML Best Paper Honourable Mention"*

Our key contributions are:  
To provide the *first rigorous proof* of the linear relationship between  
word embeddings of analogies ...

"ICML Best Paper Honourable Mention"

# Pre-requisites

**-Word2Vec**



<Word2Vec>

The quick brown fox jumps over the lazy dog

the	1	0	0	0	0	0	0	0
quick	0	1	0	0	0	0	0	0
brown	0	0	1	0	0	0	0	0
fox	0	0	0	1	0	0	0	0
jumps	0	0	0	0	1	0	0	0
over	0	0	0	0	0	1	0	0
lazy	0	0	0	0	0	0	1	0
dog	0	0	0	0	0	0	0	1

## Pre-requisites

-Word2Vec

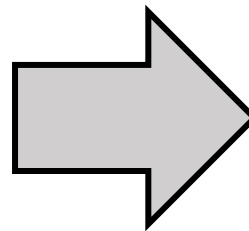
### <Word2Vec>



(the, quick)

(the, brown)

<Training Samples>



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

<Target Word>

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

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

<Context Word>

## Pre-requisites

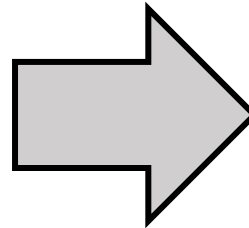
-Word2Vec

### <Word2Vec>



(quick, the)  
(quick, brown)  
(quick, fox)

<Training Samples>



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

<Target Word>

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

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

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

<Context Word>

## Pre-requisites

-Word2Vec

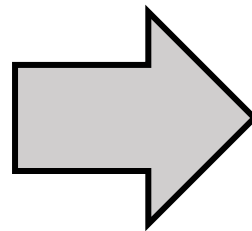
### <Word2Vec>



(dog, the)

(dog, lazy)

<Training Samples>



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

<Target Word>

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

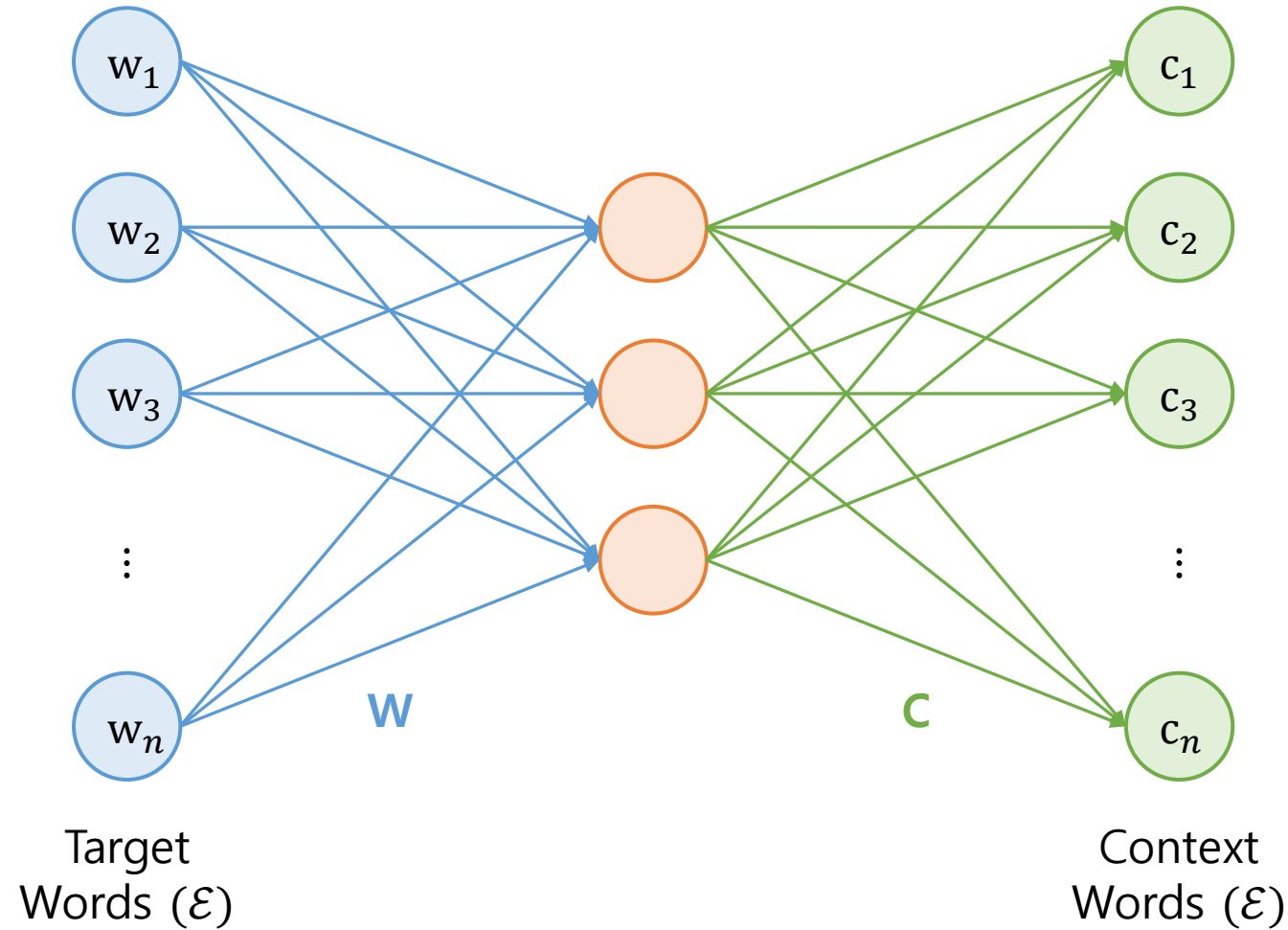
[0, 0, 0, 0, 0, 0, 1, 0]

<Context Word>

## Pre-requisites

-Word2Vec

### <Word2Vec>



## Pre-requisites

-Word2Vec

<Word2Vec>

0	0	0	1	0	0	0	0
---	---	---	---	---	---	---	---

<One-hot Encoding of "fox">

×

0.1	0.4	0.3	1.5
0.1	0.2	1.6	0.3
0.8	0.3	0.4	1.4
1.4	0.4	0.6	0.9
0.6	1.4	0.9	1.4
0.6	0.2	0.5	1.6
1.0	1.4	1.4	0.9
1.5	1.2	0.7	0.7

<Matrix  $W^T$ >

=

1.4	0.4	0.6	0.9
-----	-----	-----	-----

<Embedding of "fox">

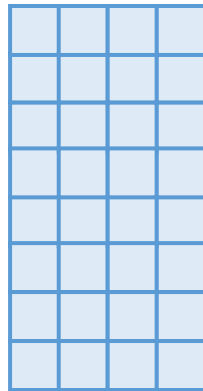
## &lt;Word Embedding as Matrix Factorization&gt;

$$W^T C \approx \text{PMI} - \log k = \text{SPMI}$$

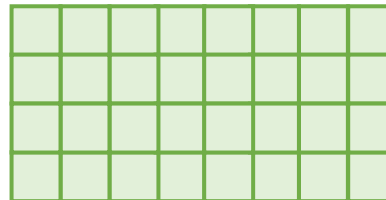
$$\text{PMI}_{i,j} = \log \frac{p(w_i, c_j)}{p(w_i)p(c_j)}$$

$w_i, c_j$  : column of  $W, C$

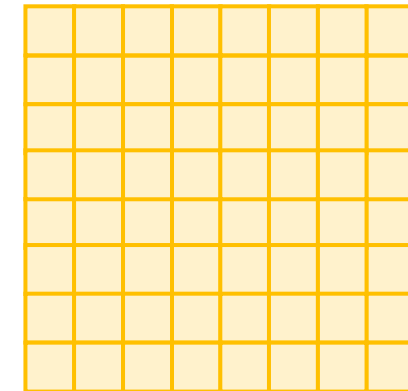
where,  $k$  is chosen number of negative samples

<Matrix  $W^T$ >

×

<Matrix  $C$ >

≈



&lt;SPMI&gt;

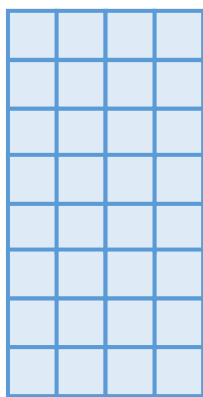
## Pre-requisites

-Word2Vec

### <Word Embedding as Matrix Factorization>

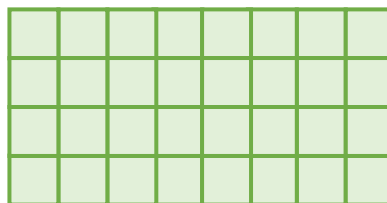
$$W^T C \approx \text{PMI} - \log k = \text{SPMI}$$

$$\text{PMI}_{i,j} = \log \frac{p(w_i, c_j)}{p(w_i)p(c_j)}$$



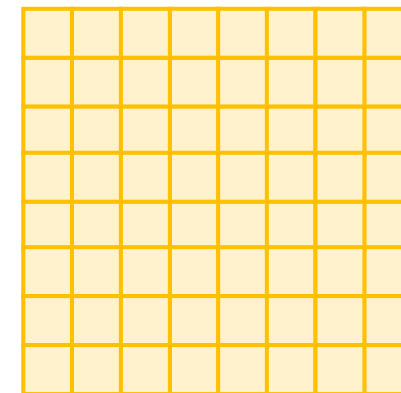
<Matrix  $W^T$ >

×



<Matrix  $C$ >

≈



<SPMI>

"We analyze skip-gram with negative-samples and show that it is implicitly factorizing a word-context matrix, whose cells are the pointwise mutual information (PMI) of the respective word and context pairs, shifted by a global constant."  
(Levy and Goldberg, 2014)



# Preliminaries

- Impact of the Shift
- Reconstruction Error
- Zero Co-occurrence Counts

## Preliminaries

-Impact of the Shift

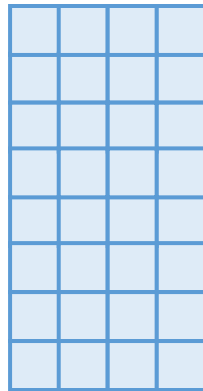
### <Impact of the Shift>

$$W^T C \approx \text{PMI} - \log k = \text{SPMI}$$

$$\text{PMI}_{i,j} = \log \frac{p(w_i, c_j)}{p(w_i)p(c_j)}$$

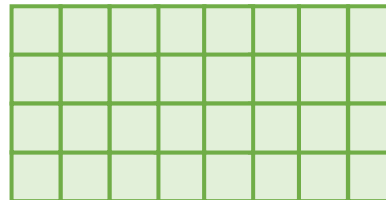
$w_i, c_j$  : column of  $W, C$

where,  $k$  is chosen number of negative samples



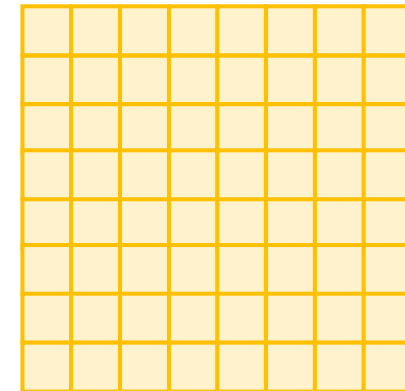
<Matrix  $W^T$ >

×



<Matrix  $C$ >

≈



<SPMI>

## Preliminaries

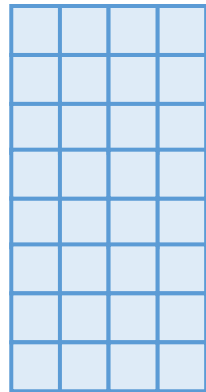
-Impact of the Shift

### <Impact of the Shift>

$$W^T C \approx \text{PMI} - \log k = \text{SPMI}$$

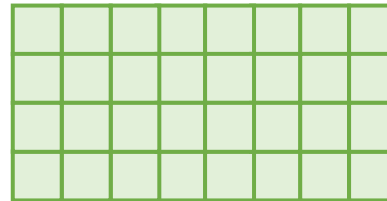
$$\text{PMI}_{i,j} = \log \frac{p(w_i, c_j)}{p(w_i)p(c_j)}$$

$w_i, c_j$  : column of  $W, C$



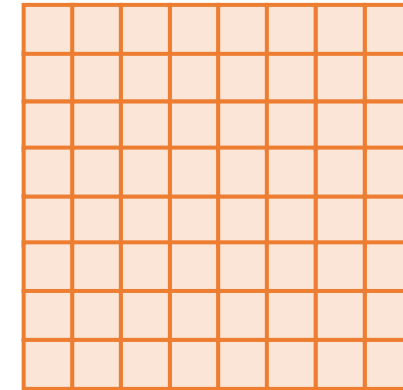
<Matrix  $W^T$ >

×



<Matrix  $C$ >

≈



<PMI>

"It is observed that adjusting the Word2Vec algorithm to avoid any direct impact of the *shift* improves embedding performance."

(Le, 2017)

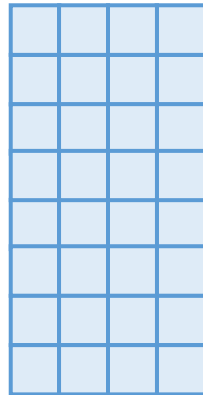
## Preliminaries

-Reconstruction Error

### <Reconstruction Error>

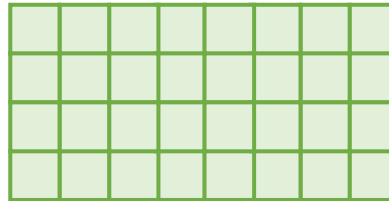
$$W^T C \approx \text{PMI}$$

$w_i, c_j$  : column of  $W, C$



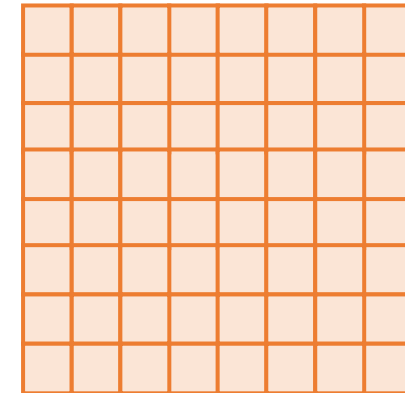
<Matrix  $W^T$ >

×



<Matrix  $C$ >

≈



<PMI>

## Preliminaries

-Reconstruction Error

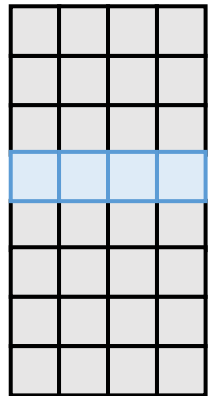
### <Reconstruction Error>

$$W^T C \approx \text{PMI}$$

$w_i, c_j$  : column of  $W, C$

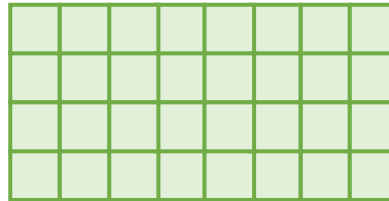
$$\mathbf{w}_i^T \mathbf{C} \approx \text{PMI}_i$$

where  $\text{PMI}_i$ : **row** of PMI



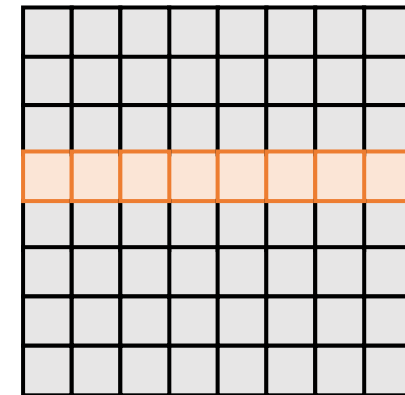
$\langle \mathbf{w}_i^T \rangle$

$\times$



$\langle \text{Matrix } C \rangle$

$\approx$



$\langle \text{PMI}_i \rangle$

## Preliminaries

-Reconstruction Error

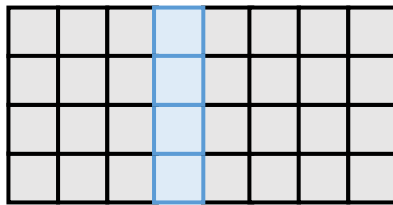
### <Reconstruction Error>

$$W^T C \approx \text{PMI}$$

$w_i, c_j$  : column of  $W, C$

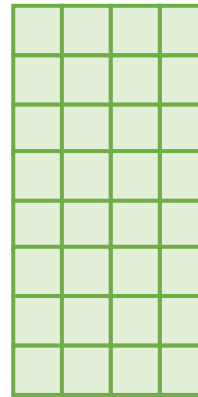
$$\mathbf{w}_i \mathbf{C}^T \approx \text{PMI}_i$$

where  $\text{PMI}_i$ : **column** of PMI



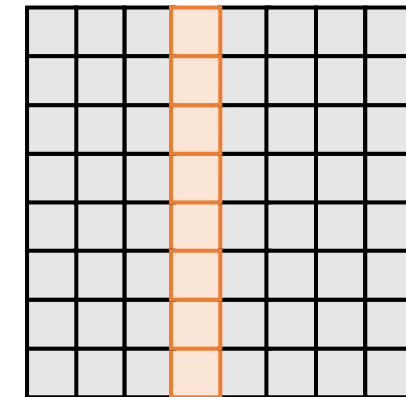
$\langle \mathbf{w}_i \rangle$

$\times$



$\langle \text{Matrix } C^T \rangle$

$\approx$



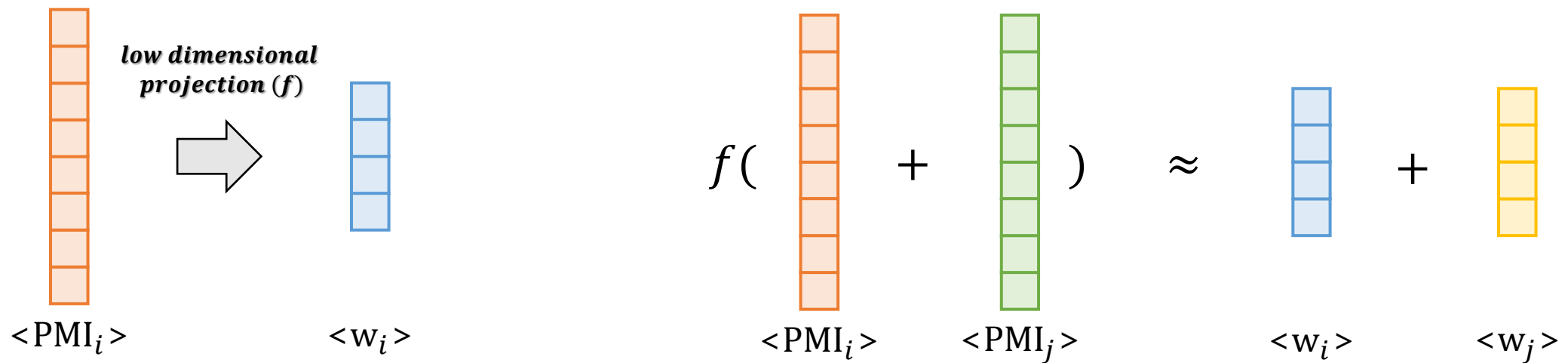
$\langle \text{PMI}_i \rangle$

## Preliminaries

-Reconstruction Error

### <Reconstruction Error>

**Assumption A2.** Letting  $\text{PMI}_k$  denote the  $k^{\text{th}}$  column of  $\text{PMI} \in \mathbb{R}^{n \times n}$ , the projection  $f: \mathbb{R}^n \rightarrow \mathbb{R}^d$ ,  $f(\text{PMI}_i) = w_i$  is approximately homomorphic with respect to addition, i.e.  $f(\text{PMI}_i + \text{PMI}_j) \approx f(\text{PMI}_i) + f(\text{PMI}_j)$



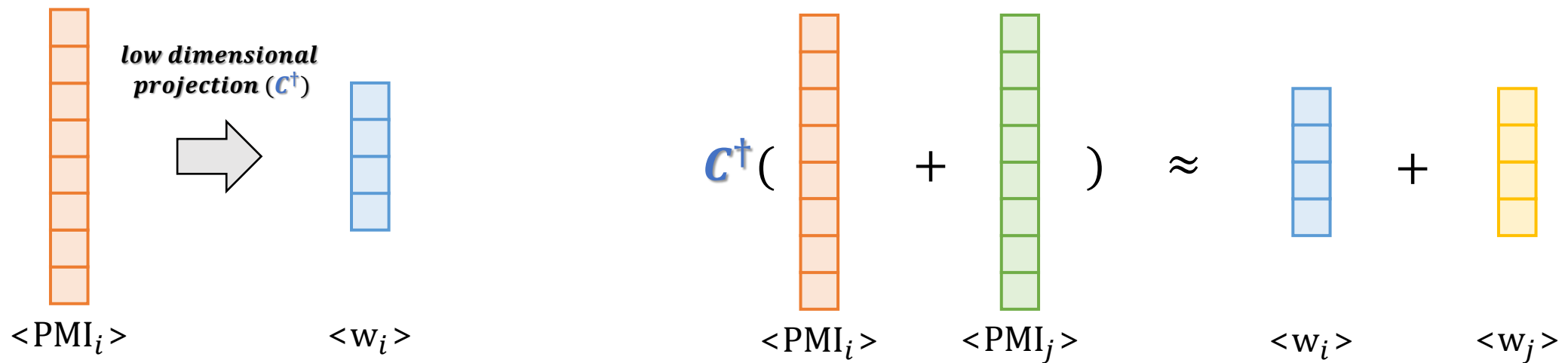
"**A2** means that whatever factorization method used, linear relationships between columns of PMI are sufficiently preserved by columns of  $W$ "

## Preliminaries

-Reconstruction Error

### <Reconstruction Error>

**Assumption A2.** Letting  $\text{PMI}_k$  denote the  $k^{\text{th}}$  column of  $\text{PMI} \in \mathbb{R}^{n \times n}$ , the projection  $\mathbf{C}^\dagger: \mathbb{R}^n \rightarrow \mathbb{R}^d$ ,  $f(\text{PMI}_i) = \mathbf{w}_i$  is approximately homomorphic with respect to addition, i. e.  $\mathbf{C}^\dagger(\text{PMI}_i + \text{PMI}_j) \approx \mathbf{C}^\dagger(\text{PMI}_i) + \mathbf{C}^\dagger(\text{PMI}_j)$



"For example, minimizing a **least squares loss function** gives linear projection  $\mathbf{w}_i = f_{LSQ}(\text{PMI}_i) = \mathbf{C}^\dagger \text{PMI}_i$   
Where  $\mathbf{C}^\dagger = (\mathbf{C}\mathbf{C}^T)^{-1}\mathbf{C}$ . We write  $f(\cdot) = \mathbf{C}^\dagger(\cdot)$  to emphasize linearity of the relationship"



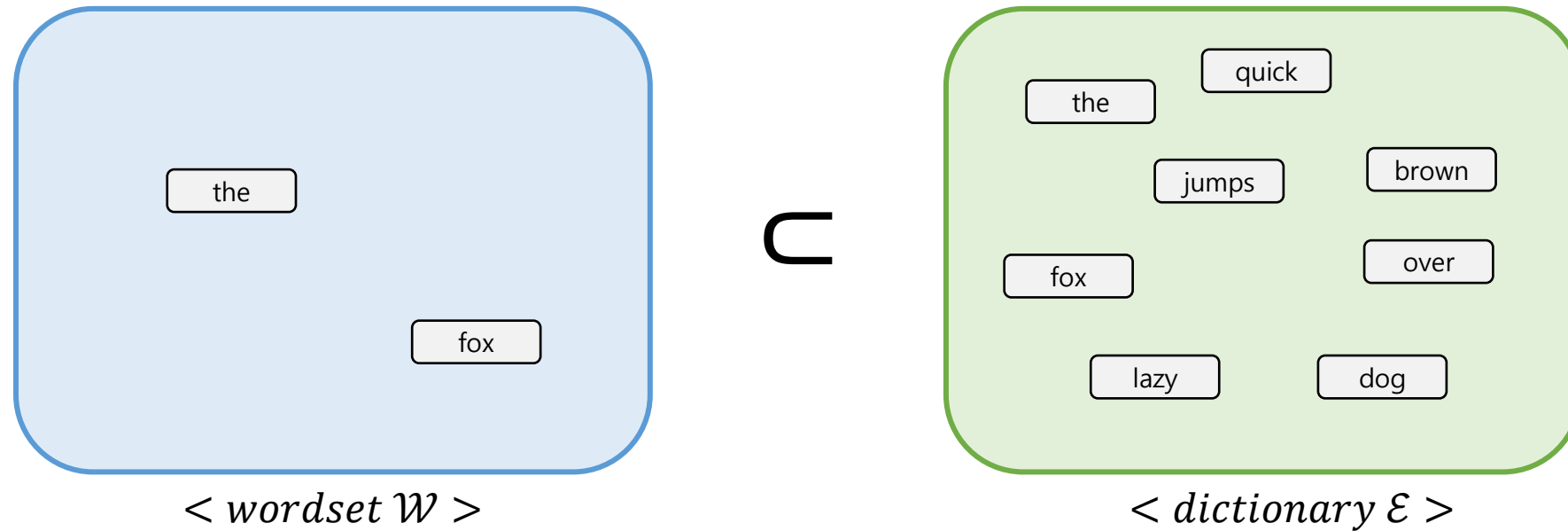
## Preliminaries

-Zero Co-occurrence Counts

### <Zero Co-occurrence counts>

$\mathcal{W}$  : small word set

$\mathcal{E}$  : fixed size dictionary



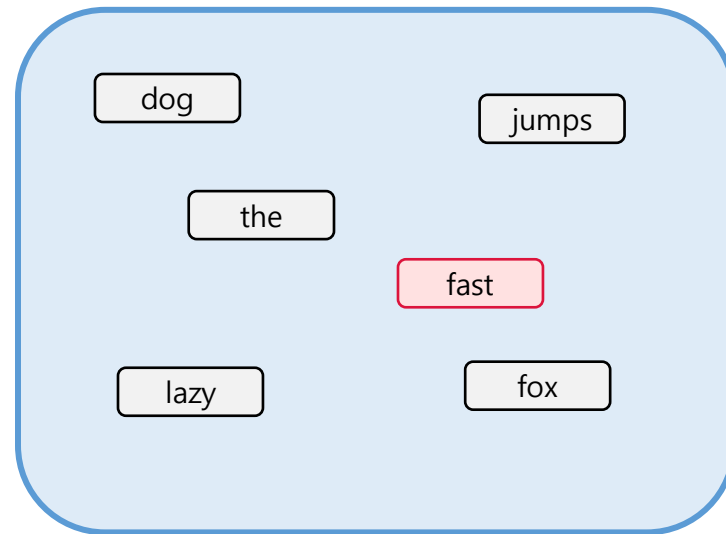
## Preliminaries

-Zero Co-occurrence Counts

### <Zero Co-occurrence counts>

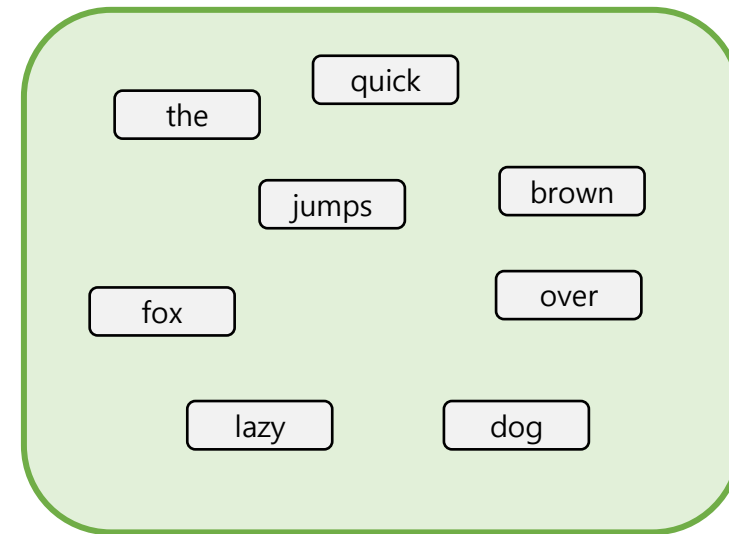
$\mathcal{W}$  : *small word set*

$\mathcal{E}$  : *fixed size dictionary*



< wordset  $\mathcal{W}$  >

$\nsubseteq$



< dictionary  $\mathcal{E}$  >

## Preliminaries

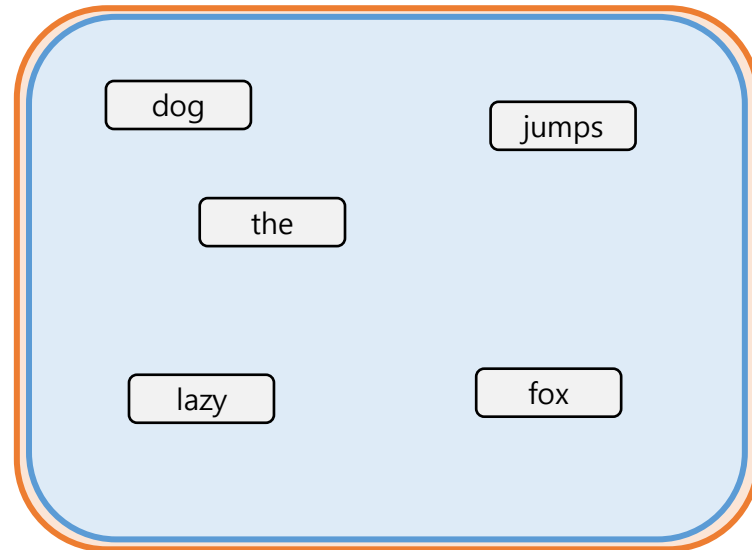
-Zero Co-occurrence Counts

### <Zero Co-occurrence counts>

$\mathcal{W}$  : small word set

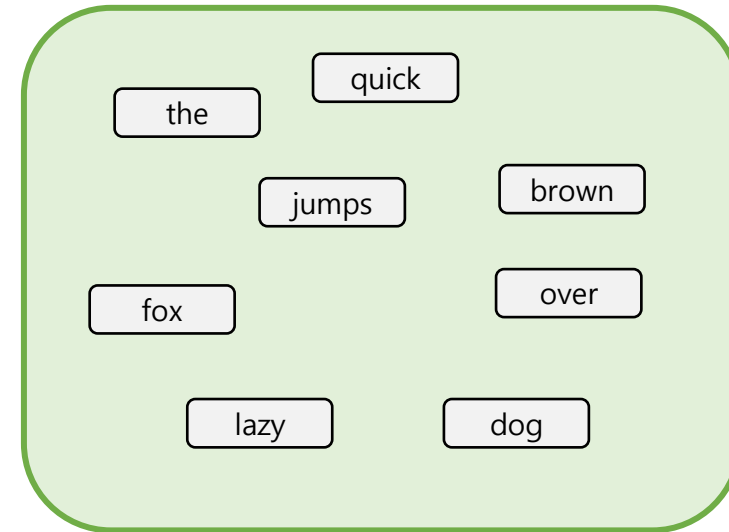
$\mathcal{E}$  : fixed size dictionary

*limit of size*



< wordset  $\mathcal{W}$  >

$\cup$



< dictionary  $\mathcal{E}$  >

## Preliminaries

-Zero Co-occurrence Counts

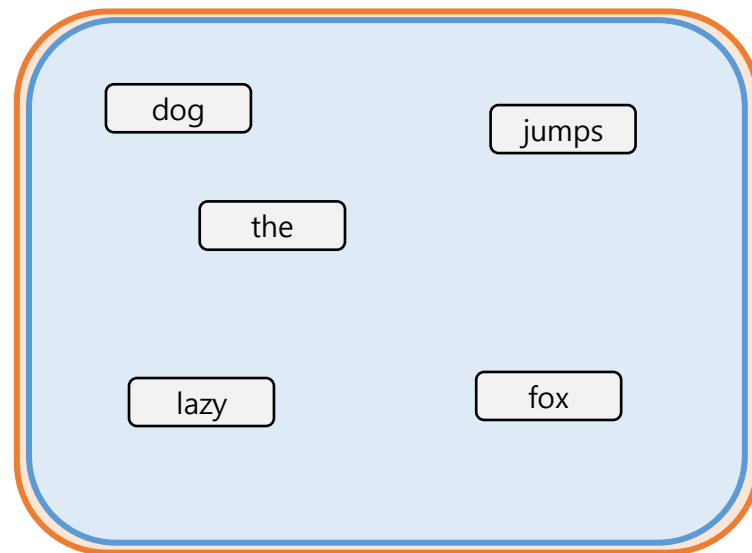
### <Zero Co-occurrence counts>

**Assumption A3.**  $p(\mathcal{W}) > 0, \forall \mathcal{W} \subseteq \mathcal{E}, |\mathcal{W}| < l$

$\mathcal{W}$  : small word set

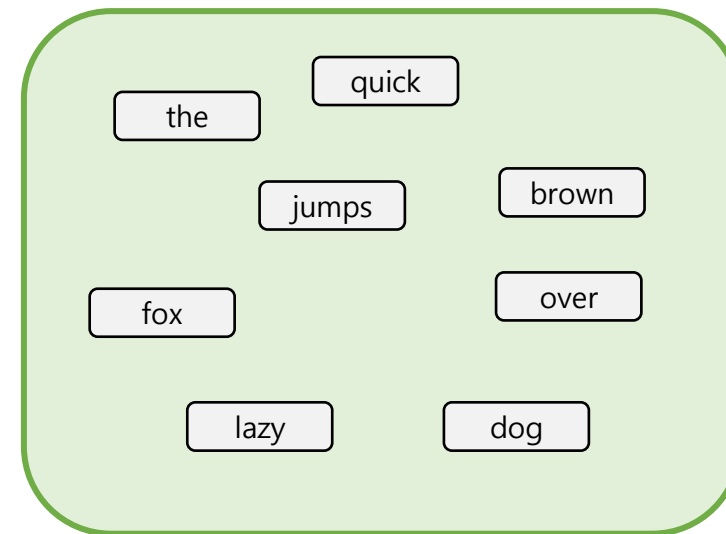
$\mathcal{E}$  : fixed size dictionary

*limit of size*



< wordset  $\mathcal{W}$  >

$\cup$



< dictionary  $\mathcal{E}$  >

# Proof

- Paraphrase
- Analogy

# Proof

- **Paraphrase**
- *Analogy*

**Proof**

-Routemap

**<Routemap>**

"man is to king as woman is to queen"



man *transforms to* king as woman *transforms to* queen



{woman, king} *paraphrases* {man, queen}



$$\text{PMI}_{\text{king}} - \text{PMI}_{\text{man}} + \text{PMI}_{\text{woman}} \approx \text{PMI}_{\text{queen}}$$



$$W_{\text{king}} - W_{\text{man}} + W_{\text{woman}} \approx W_{\text{queen}}$$

**Proof**

-Routemap

**<Routemap>**

"man is to king as woman is to queen"



man *transforms to* king as woman *transforms to* queen



{woman, king} *paraphrases* {man, queen}



$$\text{PMI}_{\text{king}} - \text{PMI}_{\text{man}} + \text{PMI}_{\text{woman}} \approx \text{PMI}_{\text{queen}}$$



$$\mathbf{w}_{\text{king}} - \mathbf{w}_{\text{man}} + \mathbf{w}_{\text{woman}} \approx \mathbf{w}_{\text{queen}}$$

semantic

geometric



**Proof**

-Routemap

**<Routemap>**

"man is to king as woman is to queen"

man *transforms to* king as woman *transforms to* queen{woman, king} *paraphrases* {man, queen}

$$\text{PMI}_{\text{king}} - \text{PMI}_{\text{man}} + \text{PMI}_{\text{woman}} \approx \text{PMI}_{\text{queen}}$$



$$\text{PMI}_i \approx \mathbf{w}_i^T \mathbf{C}$$

$$\mathbf{w}_{\text{king}} - \mathbf{w}_{\text{man}} + \mathbf{w}_{\text{woman}} \approx \mathbf{w}_{\text{queen}}$$

semantic

geometric

**Proof**

-Routemap

**<Routemap>**

"man is to king as woman is to queen"

man *transforms to* king as woman *transforms to* queen{woman, king} *paraphrases* {man, queen}

$$\text{PMI}_{\text{king}} - \text{PMI}_{\text{man}} + \text{PMI}_{\text{woman}} \approx \text{PMI}_{\text{queen}}$$



$$\text{PMI}_i \approx \mathbf{w}_i^T \mathbf{C}$$

$$\mathbf{w}_{\text{king}} - \mathbf{w}_{\text{man}} + \mathbf{w}_{\text{woman}} \approx \mathbf{w}_{\text{queen}}$$

semantic

geometric

**<Paraphrase>**

$$\mathcal{W} = \{w_1, \dots, w_m\} \subseteq \mathcal{E}$$

$w_* \in \mathcal{E}$  ***paraphrases***  $\mathcal{W}$ , if  $w_*$  and  $\mathcal{W}$  are semantically interchangeable within the text

## Proof

-Paraphrase

### <Paraphrase>

$$\mathcal{W} = \{w_1, \dots, w_m\} \subseteq \mathcal{E}$$

$w_* \in \mathcal{E}$  **paraphrases**  $\mathcal{W}$ , if  $w_*$  and  $\mathcal{W}$  are semantically interchangeable within the text

"The need for profit **may** push up prices"

<Sentence 1>

$\approx$

"The need for profit **is likely to** push up prices"

<Sentence 2>

## Proof

-Paraphrase

### <Paraphrase>

$$\mathcal{W} = \{w_1, \dots, w_m\} \subseteq \mathcal{E}$$

$w_* \in \mathcal{E}$  **paraphrases**  $\mathcal{W}$ , if  $w_*$  and  $\mathcal{W}$  are **semantically interchangeable** within the text

"The need for profit **may** push up prices"

<Sentence 1>

≈

"The need for profit **is likely to** push up prices"

<Sentence 2>

**may**

≈

**be likely to**

Semantically  
Interchangeable

## Proof

-Paraphrase

### <Paraphrase>

$$\mathcal{W} = \{w_1, \dots, w_m\} \subseteq \mathcal{E}$$

$w_* \in \mathcal{E}$  **paraphrases**  $\mathcal{W}$ , if  $w_*$  and  $\mathcal{W}$  are semantically interchangeable within the text

"The need for profit **may** push up prices"

<Sentence 1>

$\approx$

"The need for profit **is likely to** push up prices"

<Sentence 2>

**may**  
( $w_*$ )

**paraphrases**

**be likely to**  
( $\mathcal{W}$ )

## Proof

-Paraphrase

### <Paraphrase>

$$\mathcal{W} = \{w_1, \dots, w_m\} \subseteq \mathcal{E}$$

$w_* \in \mathcal{E}$  **paraphrases**  $\mathcal{W}$ , if  $w_*$  and  $\mathcal{W}$  are **semantically interchangeable** within the text

"The need for profit **may** push up prices"

<Sentence 1>

$\approx$

"The need for profit is likely to push up prices"

<Sentence 2>

may  
( $w_*$ )

paraphrases

be likely to  
( $\mathcal{W}$ )

$p(\mathbf{c}_j | w_*)$

$\approx$

$p(\mathbf{c}_j | \mathcal{W})$

## <Defining a Paraphrase>

$C_{\mathcal{W}} = \{c_{j_1}, \dots, c_{j_t}\} : \text{sequence of words (with repetition) observed in the context of } \mathcal{W}$

$w_* \in \mathcal{E} : \text{which best explains the observation of } C_{\mathcal{W}}$



## <Defining a Paraphrase>

$C_{\mathcal{W}} = \{c_{j_1}, \dots, c_{j_t}\}$  : sequence of words (*with repetition*) observed in the context of  $\mathcal{W}$

$w_* \in \mathcal{E}$  : which best explains the observation of  $C_{\mathcal{W}}$

**“The need for profit is likely to push up prices”**

## &lt;Defining a Paraphrase&gt;

$c_{\mathcal{W}} = \{c_{j_1}, \dots, c_{j_t}\}$  : sequence of words (with repetition) observed in the context of  $\mathcal{W}$

$w_* \in \mathcal{E}$  : which best explains the observation of  $c_{\mathcal{W}}$

“The need for profit is likely to push up prices”



“The need for profit may push up prices”

## &lt;Defining a Paraphrase&gt;

$C_{\mathcal{W}} = \{c_{j_1}, \dots, c_{j_t}\}$  : sequence of words (*with repetition*) observed in the context of  $\mathcal{W}$

$w_* \in \mathcal{E}$  : which best explains the observation of  $C_{\mathcal{W}}$

$$w_* = \underset{w_i \in \mathcal{E}}{\operatorname{argmax}} p(C_{\mathcal{W}} | w_i)$$

**“The need for profit may push up prices”**



## &lt;Defining a Paraphrase&gt;

$C_{\mathcal{W}} = \{c_{j_1}, \dots, c_{j_t}\}$  : sequence of words (with repetition) observed in the context of  $\mathcal{W}$

$w_* \in \mathcal{E}$  : which best explains the observation of  $C_{\mathcal{W}}$

$$w_* = \operatorname{argmax}_{w_i \in \mathcal{E}} p(C_{\mathcal{W}} | w_i)$$

assuming  $c_j \in C_{\mathcal{W}}$  to be independent draws from  $p(c_j | \mathcal{W})$

$$w_* = \operatorname{argmax}_{w_i} \prod_{c_j \in \mathcal{E}} p(c_j | w_i)^{\#j}, \#j : \text{count of } c_j \text{ in } C_{\mathcal{W}}$$

**“The need for profit may push up prices”**



## <Defining a Paraphrase>

$C_{\mathcal{W}} = \{c_{j_1}, \dots, c_{j_t}\}$  : sequence of words (with repetition) observed in the context of  $\mathcal{W}$

$w_* \in \mathcal{E}$  : which best explains the observation of  $C_{\mathcal{W}}$

assuming  $c_j \in C_{\mathcal{W}}$  to be independent draws from  $p(c_j|\mathcal{W})$

$$w_* = \operatorname{argmax}_{w_i} \prod_{c_j \in \mathcal{E}} p(c_j|w_i)^{\#j}, \#j : \text{count of } c_j \text{ in } C_{\mathcal{W}}$$

$$\rightarrow w_* = \operatorname{argmax}_{w_i} \sum_{c_j \in \mathcal{E}} p(c_j|\mathcal{W}) \log p(c_j|w_i)$$

**"The need for profit may push up prices"**



# Proof

-Paraphrase

## <Defining a Paraphrase>

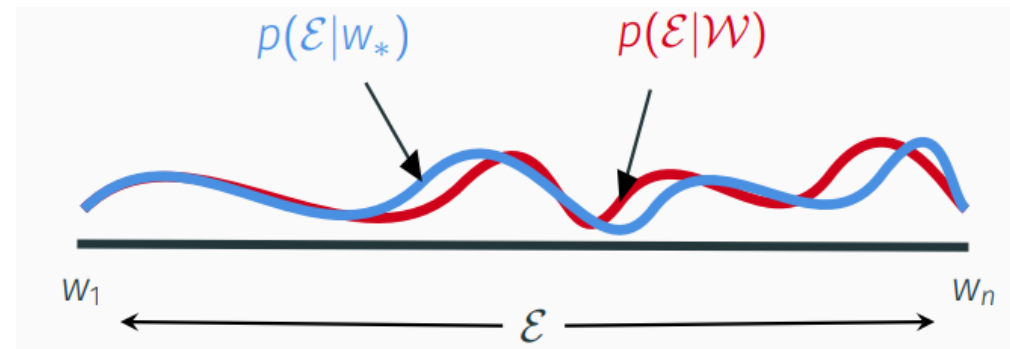
$$w_* = \operatorname{argmax}_{w_i} \prod_{c_j \in \mathcal{E}} p(c_j | w_i)^{\#j}$$

$$\rightarrow w_* = \operatorname{argmax}_{w_i} \sum_{c_j \in \mathcal{E}} p(c_j | \mathcal{W}) \log p(c_j | w_i)$$

Minimize  
KL-Divergence

$$\Delta_{KL}^{\mathcal{W}, w_*} = D_{KL}[P(c_j | \mathcal{W}) || P(c_j | w_*)]$$

$$= \sum_j p(c_j | \mathcal{W}) \log \frac{p(c_j | \mathcal{W})}{p(c_j | w_*)}$$



"the KL divergence lower bound (zero) is achieved *iff* the induced distributions are equal"

$$\Delta_{KL}^{\mathcal{W}, w_*} = 0 \leftrightarrow p(c_j | w_*) = p(c_j | \mathcal{W})$$

## &lt;Defining a Paraphrase&gt;

**Definition D1.** We say word  $w_* \in \mathcal{E}$  **paraphrases** word set  $\mathcal{W} \subseteq \mathcal{E}, |\mathcal{W}| < l$ , if the **paraphrase error**  $\rho^{\mathcal{W}, w_*} \in \mathbb{R}^n$  is (element-wise) small, where:

$$\rho_j^{\mathcal{W}, w_*} = \log \frac{p(c_j | w_*)}{p(c_j | \mathcal{W})}, c_j \in \mathcal{E}$$

"The need for profit **may** push up prices"

<Sentence 1>

"The need for profit **is likely to** push up prices"

<Sentence 2>

$\approx$

**may**  
( $w_*$ )

**paraphrases**

**be likely to**  
( $\mathcal{W}$ )

$w_* \approx_P \mathcal{W}$

## Proof

-Paraphrase

### <Paraphrase = Embedding Sum + Error>

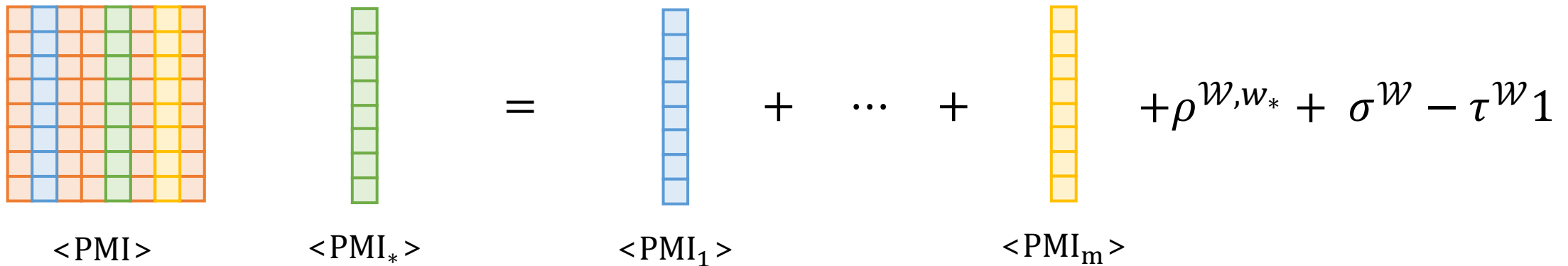
**Lemma 1.** for any word  $w_* \in \mathcal{E}$  and word set  $\mathcal{W} \subseteq \mathcal{E}, |\mathcal{W}| < l$ :

$$\text{PMI}_* = \sum_{w_i \in \mathcal{W}} \text{PMI}_i + \rho^{\mathcal{W}, w_*} + \sigma^{\mathcal{W}} - \tau^{\mathcal{W}} \mathbf{1}$$

$\text{PMI}_\bullet$  : column of PMI corresponding to  $w_\bullet \in \mathcal{E}$

$\mathbf{1} \in \mathbb{R}^n$  : vector of 1s

$$\sigma_j^{\mathcal{W}} = \log \frac{p(\mathcal{W} | c_j)}{\prod_i p(w_i | c_j)}, \tau^{\mathcal{W}} = \log \frac{p(\mathcal{W})}{\prod_i p(w_i)} : \text{error terms}$$



<PMI>      <PMI<sub>\*</sub>>      =      <PMI<sub>1</sub>>      +      ...      +      <PMI<sub>m</sub>>      +  $\rho^{\mathcal{W}, w_*} + \sigma^{\mathcal{W}} - \tau^{\mathcal{W}} \mathbf{1}$



<Paraphrase = Embedding Sum + Error>

*Proof of Lemma 1.*

$$\begin{aligned}
 & PMI(w_*, c_j) - \sum_{w_i \in \mathcal{W}} PMI(w_i, c_j) \\
 &= \log \frac{p(w_* | c_j)}{p(w_*)} - \log \prod_{w_i \in \mathcal{W}} \frac{p(w_i | c_j)}{p(w_i)} \\
 &= \log \frac{p(w_* | c_j)}{\prod_{\mathcal{W}} p(w_i | c_j)} - \log \frac{p(w_*)}{\prod_{\mathcal{W}} p(w_i)} + \log \frac{p(\mathcal{W} | c_j)}{p(\mathcal{W} | c_j)} + \log \frac{p(\mathcal{W})}{p(\mathcal{W})} \\
 &= \log \frac{p(w_* | c_j)}{p(\mathcal{W} | c_j)} - \log \frac{p(w_*)}{p(\mathcal{W})} + \log \frac{p(\mathcal{W} | c_j)}{\prod_{\mathcal{W}} p(w_i | c_j)} - \log \frac{p(\mathcal{W})}{\prod_{\mathcal{W}} p(w_i)} \\
 &= \log \frac{p(c_j | w_*)}{p(c_j | \mathcal{W})} + \log \frac{p(\mathcal{W} | c_j)}{\prod_{\mathcal{W}} p(w_i | c_j)} - \log \frac{p(\mathcal{W})}{\prod_{\mathcal{W}} p(w_i)} \\
 &= \rho_j^{\mathcal{W}, w_*} + \sigma_j^{\mathcal{W}} - \tau^{\mathcal{W}}
 \end{aligned}$$

## &lt;Paraphrase = Embedding Sum + Error&gt;

**Lemma 1.** for any word  $w_* \in \mathcal{E}$  and word set  $\mathcal{W} \subseteq \mathcal{E}, |\mathcal{W}| < l$ :

$$\text{PMI}_* = \sum_{w_i \in \mathcal{W}} \text{PMI}_i + \rho^{\mathcal{W}, w_*} + \sigma^{\mathcal{W}} - \tau^{\mathcal{W}} 1$$



$$C^\dagger(\text{PMI}_*) = C^\dagger \left( \sum_{w_i \in \mathcal{W}} \text{PMI}_i + \rho^{\mathcal{W}, w_*} + \sigma^{\mathcal{W}} - \tau^{\mathcal{W}} 1 \right)$$

$$w_* = w_{\mathcal{W}} + C^\dagger(\rho^{\mathcal{W}, w_*} + \sigma^{\mathcal{W}} - \tau^{\mathcal{W}} 1)$$

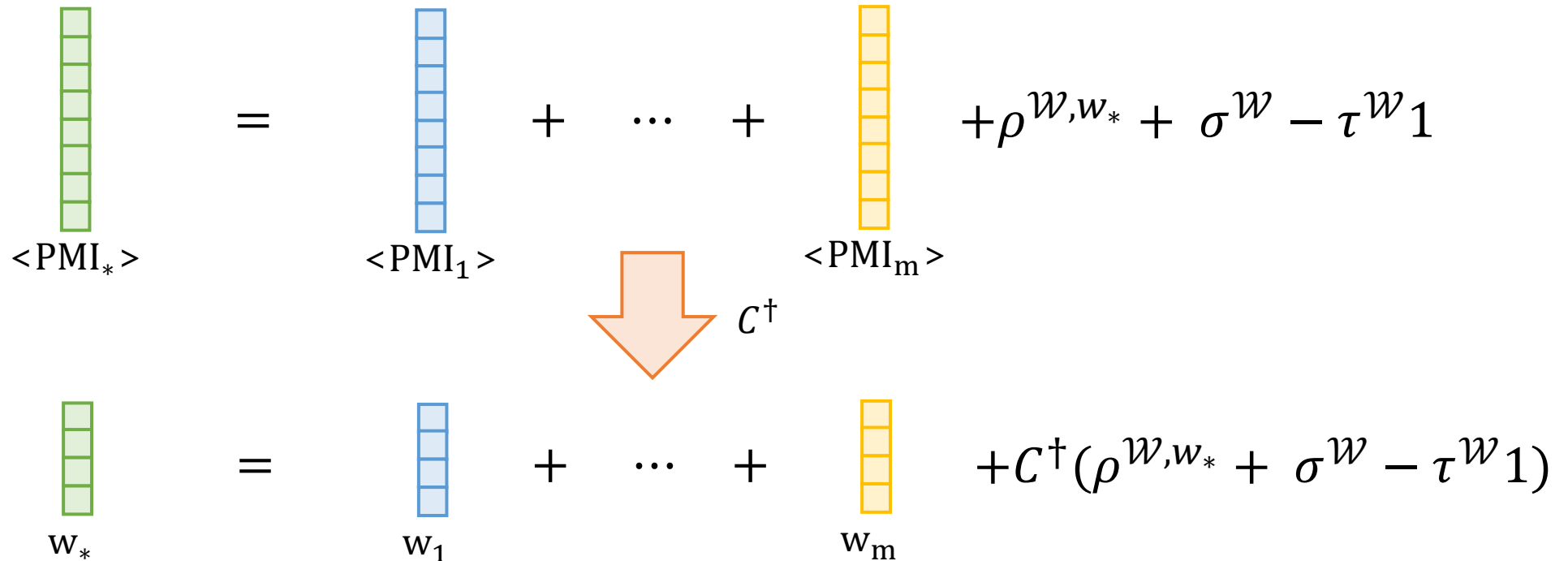
$$\text{where, } w_* = \sum_{w_i \in \mathcal{W}} w_i$$

## <Paraphrase = Embedding Sum + Error>

**Theorem 1.** (Paraphrase). for any word  $w_* \in \mathcal{E}$  and word set  $\mathcal{W} \subseteq \mathcal{E}, |\mathcal{W}| < l$ :

$$w_* = w_{\mathcal{W}} + C^\dagger(\rho^{\mathcal{W}, w_*} + \sigma^{\mathcal{W}} - \tau^{\mathcal{W}} 1)$$

$$\text{where, } w_{\mathcal{W}} = \sum_{w_i \in \mathcal{W}} w_i$$

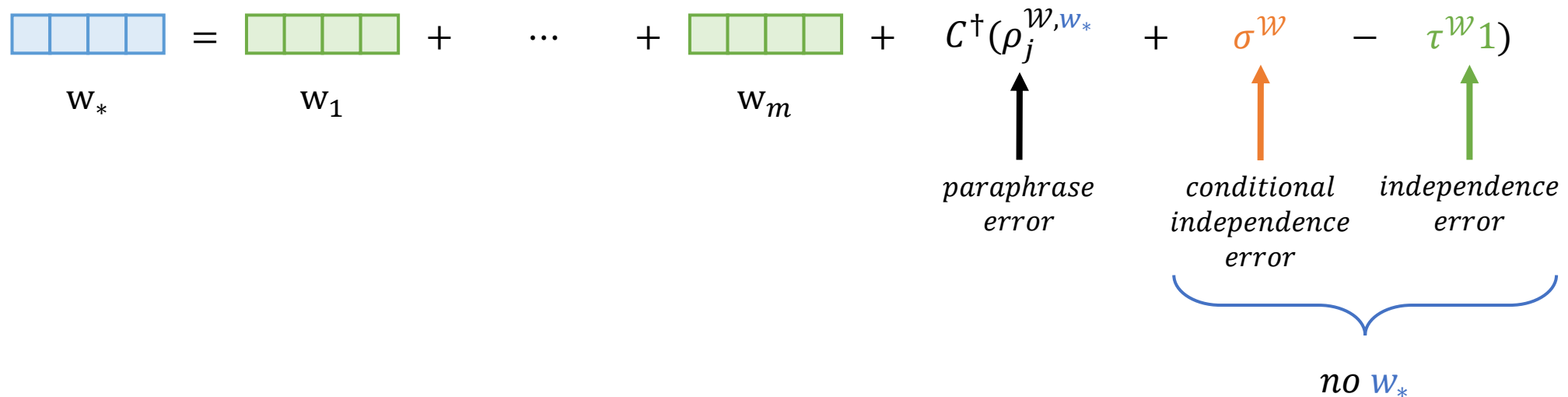


# <Paraphrase = Embedding Sum + Error>

**Theorem 1.** (Paraphrase). for any word  $w_* \in \mathcal{E}$  and word set  $\mathcal{W} \subseteq \mathcal{E}, |\mathcal{W}| < l$ :

$$w_* = w_{\mathcal{W}} + C^\dagger(\rho^{\mathcal{W}, w_*} + \sigma^{\mathcal{W}} - \tau^{\mathcal{W}} 1)$$

$$\text{where, } w_{\mathcal{W}} = \sum_{w_i \in \mathcal{W}} w_i$$



## <Paraphrase = Embedding Sum + Error>

**Theorem 1.** (Paraphrase). for any word  $w_* \in \mathcal{E}$  and word set  $\mathcal{W} \subseteq \mathcal{E}, |\mathcal{W}| < l$ :

$$w_* = w_{\mathcal{W}} + C^\dagger(\rho^{\mathcal{W}, w_*} + \sigma^{\mathcal{W}} - \tau^{\mathcal{W}} 1)$$

$$\text{where, } w_{\mathcal{W}} = \sum_{w_i \in \mathcal{W}} w_i$$

$$\sigma_j^{\mathcal{W}} = \log \frac{p(\mathcal{W}|c_j)}{\prod_i p(w_i|c_j)} : \text{conditional independence error}$$



$p(\mathcal{W}|c_j)$

**"The need for profit is likely to push up prices"**



$\prod_i p(w_i|c_j)$

$\sigma_j^{\mathcal{W}} = 0$  iff all  $w_i \in \mathcal{W}$  are conditionally independent given each  $c_j \in \mathcal{E}$

## <Paraphrase = Embedding Sum + Error>

**Theorem 1.** (Paraphrase). for any word  $w_* \in \mathcal{E}$  and word set  $\mathcal{W} \subseteq \mathcal{E}, |\mathcal{W}| < l$ :

$$w_* = w_{\mathcal{W}} + C^\dagger(\rho^{\mathcal{W}, w_*} + \sigma^{\mathcal{W}} - \tau^{\mathcal{W}} \mathbf{1})$$

$$\text{where, } w_{\mathcal{W}} = \sum_{w_i \in \mathcal{W}} w_i$$

$$\tau^{\mathcal{W}} = \log \frac{p(\mathcal{W})}{\prod_i p(w_i)} : \text{independence error}$$



$p(\mathcal{W})$

"The need for profit **is likely to** push up prices"



$\prod_i p(w_i)$

$\tau^{\mathcal{W}} = 0$  iff all  $w_i \in \mathcal{W}$  are mutually independent

## <Paraphrase = Embedding Sum + Error>

**Theorem 1.** (Paraphrase). for any word  $w_* \in \mathcal{E}$  and word set  $\mathcal{W} \subseteq \mathcal{E}, |\mathcal{W}| < l$ :

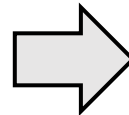
$$w_* = w_{\mathcal{W}} + C^\dagger(\rho^{\mathcal{W}, w_*} + \sigma^{\mathcal{W}} - \tau^{\mathcal{W}} \mathbf{1})$$

$$\text{where, } w_{\mathcal{W}} = \sum_{w_i \in \mathcal{W}} w_i$$

**Definition D1.** We say word  $w_* \in \mathcal{E}$  **paraphrases** word set  $\mathcal{W} \subseteq \mathcal{E}, |\mathcal{W}| < l$ , if the **paraphrase error**  $\rho^{\mathcal{W}, w_*} \in \mathbb{R}^n$  is (element-wise) **small**, where:

$$\rho_j^{\mathcal{W}, w_*} = \log \frac{p(c_j | w_*)}{p(c_j | \mathcal{W})}, c_j \in \mathcal{E}$$

king paraphrase {man, royal}



$$w_{\text{man}} + w_{\text{royal}} \approx w_{\text{king}}$$

$\rho, \sigma, \tau$

## <What We Want to Prove>

{woman, king} *paraphrases* {man, queen}



$$W_{\text{king}} - W_{\text{man}} + W_{\text{woman}} \approx W_{\text{queen}}$$

## <What We Have Proven>

king *paraphrases* {man, royal}



*Dependency  
Error*

$$W_{\text{man}} + W_{\text{royal}} \approx W_{\text{king}}$$



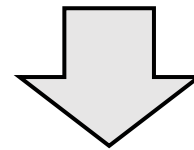
**Proof**

-Analogy

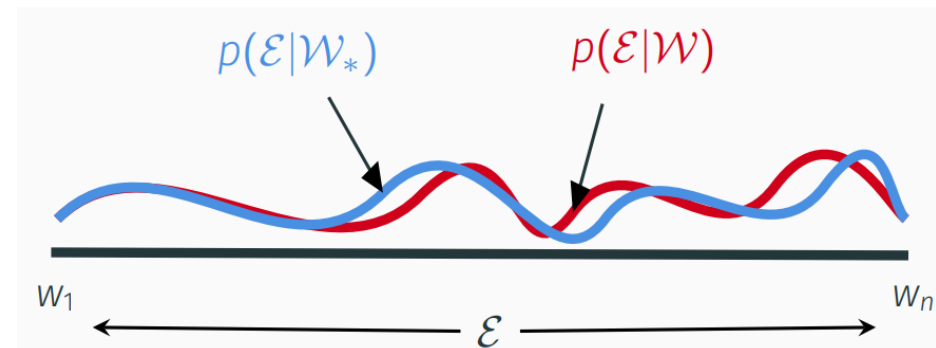
**<Paraphrasing Word Sets>**

**Definition D1.** We say word  $w_* \in \mathcal{E}$  **paraphrases** word set  $\mathcal{W} \subseteq \mathcal{E}, |\mathcal{W}| < l$ , if the **paraphrase error**  $\rho^{\mathcal{W}, w_*} \in \mathbb{R}^n$  is (element-wise) small, where:

$$\rho_j^{\mathcal{W}, w_*} = \log \frac{p(c_j | w_*)}{p(c_j | \mathcal{W})}, c_j \in \mathcal{E}$$

one word  $w_*$ word set  $\mathcal{W}_*$ 

$$\rho_j^{\mathcal{W}, \mathcal{W}_*} = \log \frac{p(c_j | \mathcal{W}_*)}{p(c_j | \mathcal{W})}, c_j \in \mathcal{E}$$



## <Paraphrasing Word Sets>

**Definition D2.** We say word set  $\mathcal{W}_* \subseteq \mathcal{E}$  **paraphrases** word set  $\mathcal{W} \subseteq \mathcal{E}$ ,  $|\mathcal{W}|, |\mathcal{W}_*| < l$ , if the **paraphrase error**  $\rho^{\mathcal{W}, \mathcal{W}_*} \in \mathbb{R}^n$  is (element-wise) small, where:

$$\rho_j^{\mathcal{W}, \mathcal{W}_*} = \log \frac{p(c_j | \mathcal{W}_*)}{p(c_j | \mathcal{W})}, c_j \in \mathcal{E}$$

"The need for profit **will probably** push up prices"

<Sentence 1>

"The need for profit **is likely to** push up prices"

<Sentence 2>

$\approx$

**will probably**  
( $\mathcal{W}_*$ )

**paraphrases**

**be likely to**  
( $\mathcal{W}$ )

$\mathcal{W}_* \approx_P \mathcal{W}$

**Proof**

-Analogy

**<Paraphrasing Word Sets>****Lemma 2.** for any word sets  $\mathcal{W}$ ,  $\mathcal{W}_* \subseteq \mathcal{E}$ ,  $|\mathcal{W}|$ ,  $|\mathcal{W}_*| < l$ :

$$\sum_{w_i \in \mathcal{W}_*} \text{PMI}_i = \sum_{w_i \in \mathcal{W}} \text{PMI}_i + \rho^{\mathcal{W}, \mathcal{W}_*} + \sigma^{\mathcal{W}} - \sigma^{\mathcal{W}_*} - (\tau^{\mathcal{W}} - \tau^{\mathcal{W}_*})\mathbf{1}$$

 $\text{PMI}_\bullet$  : column of PMI corresponding to  $w_\bullet \in \mathcal{E}$  $\mathbf{1} \in \mathbb{R}^n$  : vector of 1s

$$\sigma_j^{\mathcal{W}} = \log \frac{p(\mathcal{W}|c_j)}{\prod_i p(w_i|c_j)}, \tau^{\mathcal{W}} = \log \frac{p(\mathcal{W})}{\prod_i p(w_i)} : \text{error terms}$$

## <Paraphrasing Word Sets>

**Theorem 2.** (Generalised Paraphrase). for any word sets  $\mathcal{W}, \mathcal{W}_* \subseteq \mathcal{E}, |\mathcal{W}|, |\mathcal{W}_*| < l$ :

$$w_{\mathcal{W}_*} = w_{\mathcal{W}} + C^\dagger(\rho^{\mathcal{W}, \mathcal{W}_*} + \sigma^{\mathcal{W}} - \sigma^{\mathcal{W}_*} - (\tau^{\mathcal{W}} - \tau^{\mathcal{W}_*})1)$$

$$\text{where, } w_{\mathcal{W}} = \sum_{w_i \in \mathcal{W}} w_i$$

**Proof of Theorem 2.** Multiply **Lemma 2.** by  $C^\dagger$

## <Paraphrasing Word Sets>

**Definition D2.** We say word set  $\mathcal{W}_* \subseteq \mathcal{E}$  **paraphrases** word set  $\mathcal{W} \subseteq \mathcal{E}$ ,  $|\mathcal{W}|, |\mathcal{W}_*| < l$ , if the **paraphrase error**  $\rho^{\mathcal{W}, \mathcal{W}_*} \in \mathbb{R}^n$  is **(element-wise) small**, where:

$$\rho_j^{\mathcal{W}, \mathcal{W}_*} = \log \frac{p(c_j | \mathcal{W}_*)}{p(c_j | \mathcal{W})}, c_j \in \mathcal{E}$$

**Lemma 2.** for any word sets  $\mathcal{W}, \mathcal{W}_* \subseteq \mathcal{E}$ ,  $|\mathcal{W}|, |\mathcal{W}_*| < l$ :

$$\sum_{w_i \in \mathcal{W}_*} \text{PMI}_i = \sum_{w_i \in \mathcal{W}} \text{PMI}_i + \rho^{\mathcal{W}, \mathcal{W}_*} + \sigma^{\mathcal{W}} - \sigma^{\mathcal{W}_*} - (\tau^{\mathcal{W}} - \tau^{\mathcal{W}_*})\mathbf{1}$$

**Theorem 2.** (Generalised Paraphrase). for any word sets  $\mathcal{W}, \mathcal{W}_* \subseteq \mathcal{E}$ ,  $|\mathcal{W}|, |\mathcal{W}_*| < l$ :

$$\mathbf{w}_{\mathcal{W}_*} = \mathbf{w}_{\mathcal{W}} + \mathcal{C}^\dagger(\rho^{\mathcal{W}, \mathcal{W}_*} + \sigma^{\mathcal{W}} - \sigma^{\mathcal{W}_*} - (\tau^{\mathcal{W}} - \tau^{\mathcal{W}_*})\mathbf{1})$$



{woman, king} *paraphrases* {man, queen}

$\Downarrow (\rho, \sigma, \tau)$

$$\text{PMI}_{\text{king}} - \text{PMI}_{\text{man}} + \text{PMI}_{\text{woman}} \approx \text{PMI}_{\text{queen}}$$

$\Downarrow$

$$\mathbf{w}_{\text{king}} - \mathbf{w}_{\text{man}} + \mathbf{w}_{\text{woman}} \approx \mathbf{w}_{\text{queen}}$$

**Proof**

-Routemap

**<Routemap>**

"man is to king as woman is to queen"

man *transforms to* king as woman *transforms to* queen{woman, king} *paraphrases* {man, queen}**Dependency  
Error**

$$\text{PMI}_{\text{king}} - \text{PMI}_{\text{man}} + \text{PMI}_{\text{woman}} \approx \text{PMI}_{\text{queen}}$$



$$\text{PMI}_i \approx \mathbf{w}_i^T \mathbf{C}$$

$$\mathbf{w}_{\text{king}} - \mathbf{w}_{\text{man}} + \mathbf{w}_{\text{woman}} \approx \mathbf{w}_{\text{queen}}$$

semantic

geometric

**Proof**

-Routemap

**<Routemap>**

"man is to king as woman is to queen"



man *transforms to* king as woman *transforms to* queen



{woman, king} *paraphrases* {man, queen}



**Dependency  
Error**

$$\text{PMI}_{\text{king}} - \text{PMI}_{\text{man}} + \text{PMI}_{\text{woman}} \approx \text{PMI}_{\text{queen}}$$



$$\text{PMI}_i \approx \mathbf{w}_i^T \mathbf{C}$$

$$\mathbf{w}_{\text{king}} - \mathbf{w}_{\text{man}} + \mathbf{w}_{\text{woman}} \approx \mathbf{w}_{\text{queen}}$$

semantic

geometric

# Proof

- Paraphrase
- **Analogy**



## <Analogy>

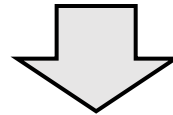
***analogy***:  $w_a$  is to  $w_{a^*}$  as  $w_b$  is to  $w_{b^*}$

where,  $w_a, w_{a^*}, w_b, w_{b^*} \in \mathcal{E}$

## <Analogy>

**analogy:**  $w_a$  is to  $w_{a^*}$  as  $w_b$  is to  $w_{b^*}$

where,  $w_a, w_{a^*}, w_b, w_{b^*} \in \mathcal{E}$



More  
Rigorously

**analogy:**  $\mathfrak{A}$

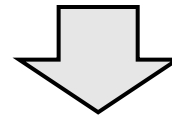
$S_{\mathfrak{A}} \subseteq \mathcal{E} \times \mathcal{E}$  : set of ordered word pairs

$(w_x, w_{x^*}) \in S_{\mathfrak{A}}$  iff " $w_x$  is to  $w_{x^*}$  as [all other analogical pairs] under  $\mathfrak{A}$ "

## &lt;Analogy&gt;

**analogy:**  $w_a$  is to  $w_{a^*}$  as  $w_b$  is to  $w_{b^*}$

where,  $w_a, w_{a^*}, w_b, w_{b^*} \in \mathcal{E}$

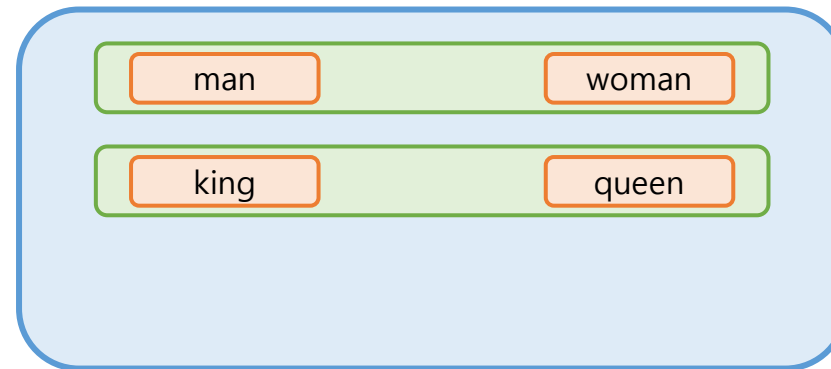


More  
Rigorously

**analogy:**  $\mathfrak{A}$

$S_{\mathfrak{A}} \subseteq \mathcal{E} \times \mathcal{E}$  : set of ordered word pairs

$(w_x, w_{x^*}) \in S_{\mathfrak{A}}$  iff " $w_x$  is to  $w_{x^*}$  as [all other analogical pairs] under  $\mathfrak{A}$ "

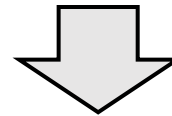


$\mathfrak{A}$

## &lt;Analogy&gt;

**analogy:**  $w_a$  is to  $w_{a^*}$  as  $w_b$  is to  $w_{b^*}$

where,  $w_a, w_{a^*}, w_b, w_{b^*} \in \mathcal{E}$

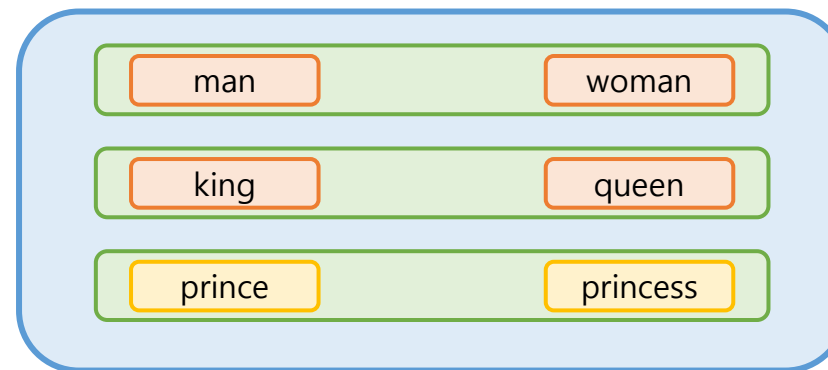


More  
Rigorously

**analogy:**  $\mathfrak{A}$

$S_{\mathfrak{A}} \subseteq \mathcal{E} \times \mathcal{E}$  : set of ordered word pairs

$(w_x, w_{x^*}) \in S_{\mathfrak{A}}$  iff " $w_x$  is to  $w_{x^*}$  as [all other analogical pairs] under  $\mathfrak{A}$ "

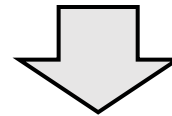


$\mathfrak{A}$

## &lt;Analogy&gt;

**analogy:**  $w_a$  is to  $w_{a^*}$  as  $w_b$  is to  $w_{b^*}$

where,  $w_a, w_{a^*}, w_b, w_{b^*} \in \mathcal{E}$

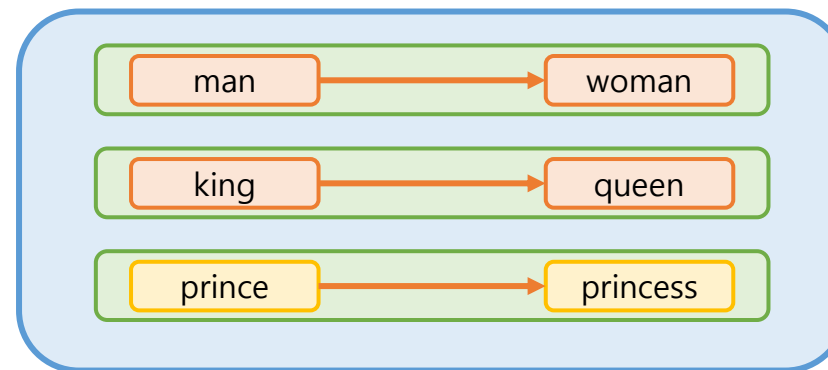


More  
Rigorously

**analogy:**  $\mathfrak{A}$

$S_{\mathfrak{A}} \subseteq \mathcal{E} \times \mathcal{E}$  : set of ordered word pairs

$(w_x, w_{x^*}) \in S_{\mathfrak{A}}$  iff " $w_x$  is to  $w_{x^*}$  as [all other analogical pairs] under  $\mathfrak{A}$ "

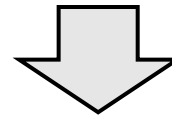


$\mathfrak{A}$

## &lt;Analogy&gt;

**analogy:**  $w_a$  is to  $w_{a^*}$  as  $w_b$  is to  $w_{b^*}$

where,  $w_a, w_{a^*}, w_b, w_{b^*} \in \mathcal{E}$

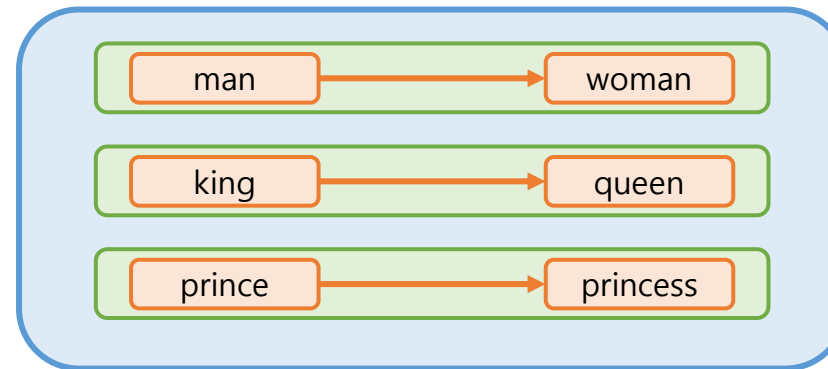


More  
Rigorously

**analogy:**  $\mathfrak{A}$

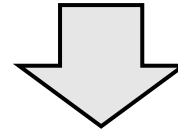
$S_{\mathfrak{A}} \subseteq \mathcal{E} \times \mathcal{E}$  : set of ordered word pairs

$(w_x, w_{x^*}) \in S_{\mathfrak{A}}$  iff " $w_x$  is to  $w_{x^*}$  as [all other analogical pairs] under  $\mathfrak{A}$ "



$\mathfrak{A}$

## &lt;Analogy&gt;

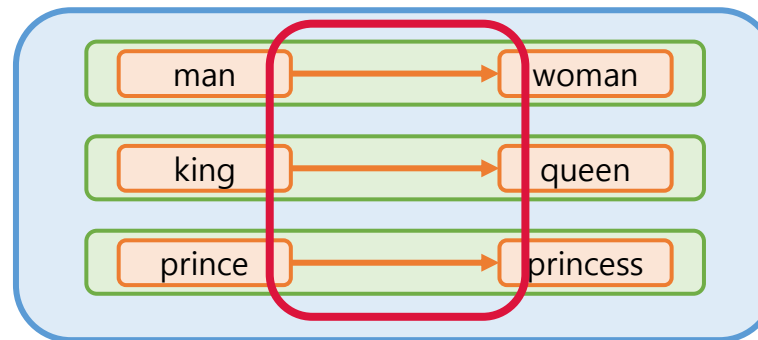
*analogy*:  $\mathfrak{A}$  $S_{\mathfrak{A}} \subseteq \mathcal{E} \times \mathcal{E}$  : set of ordered word pairs $(w_x, w_{x^*}) \in S_{\mathfrak{A}}$  iff " $w_x$  is to  $w_{x^*}$  as [all other analogical pairs] under  $\mathfrak{A}$ "

To Explain

$$w_{b^*} \approx w_{a^*} - w_a + w_b$$

In more general case

$$w_{x^*} - w_x \approx \mathbf{u}_{\mathfrak{A}}$$

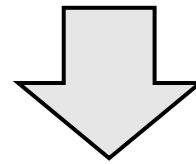
 $\forall (w_x, w_{x^*}) \in S_{\mathfrak{A}}, \mathbf{u}_{\mathfrak{A}} \in \mathbb{R}^n$ : specific vector to  $\mathfrak{A}$ 


## &lt;Word Transformation&gt;

**Theorem 2.** (Generalised Paraphrase). for any word sets  $\mathcal{W}, \mathcal{W}_* \subseteq \mathcal{E}, |\mathcal{W}|, |\mathcal{W}_*| < l$ :

$$w_{\mathcal{W}_*} = w_{\mathcal{W}} + C^\dagger(\rho^{\mathcal{W}, \mathcal{W}_*} + \sigma^{\mathcal{W}} - \sigma^{\mathcal{W}_*} - (\tau^{\mathcal{W}} - \tau^{\mathcal{W}_*})1)$$

$$\text{where, } w_{\mathcal{W}} = \sum_{w_i \in \mathcal{W}} w_i$$



$$\mathcal{W} = \{w_x\} \cup \mathcal{W}^+$$

$$\mathcal{W}_* = \{w_{x_*}\} \cup \mathcal{W}^-$$

**Corollary 2.1.** for any words  $w_x, w_{x_*} \in \mathcal{E}$  and word sets  $\mathcal{W}^+, \mathcal{W}^- \subseteq \mathcal{E}$ ,  $|\mathcal{W}^+|, |\mathcal{W}^-| < l - 1$ :

$$w_{x_*} = w_x + w_{\mathcal{W}^+} - w_{\mathcal{W}^-} + C^\dagger(\rho^{\mathcal{W}, \mathcal{W}_*} + \sigma^{\mathcal{W}} - \sigma^{\mathcal{W}_*} - (\tau^{\mathcal{W}} - \tau^{\mathcal{W}_*})1)$$

$$\text{where, } \mathcal{W} = \{w_x\} \cup \mathcal{W}^+, \mathcal{W}_* = \{w_{x_*}\} \cup \mathcal{W}^-$$



**Proof**

-Analogy

## &lt;Word Transformation&gt;

$$\mathcal{W} = \{w_x\} \cup \mathcal{W}^+, \mathcal{W}_* = \{w_{x_*}\} \cup \mathcal{W}^-$$

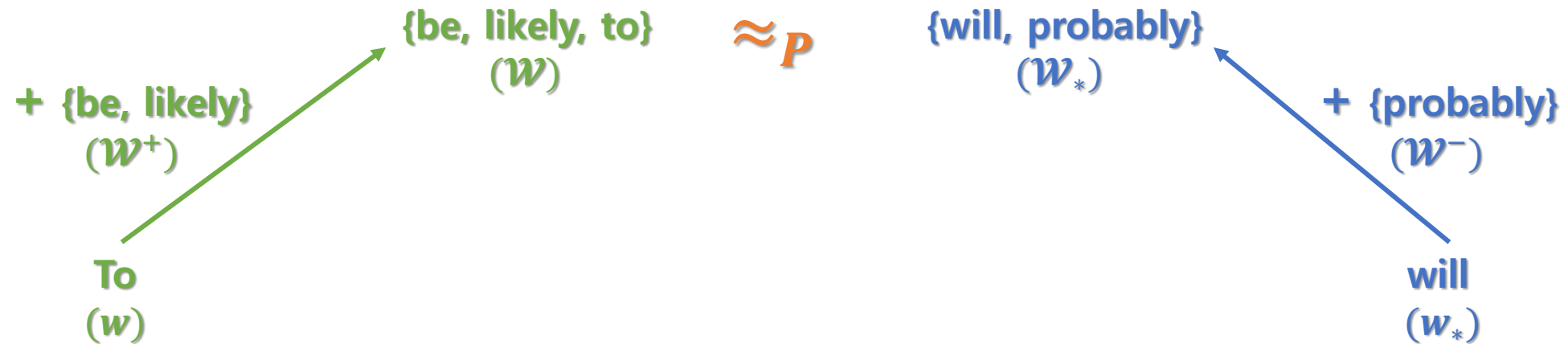
$$\begin{array}{ccc} \{\text{be, likely, to}\} & \approx_P & \{\text{will, probably}\} \\ (\mathcal{W}) & & (\mathcal{W}_*) \end{array}$$

**Proof**

-Analogy

**<Word Transformation>**

$$\mathcal{W} = \{w_x\} \cup \mathcal{W}^+, \mathcal{W}_* = \{w_{x_*}\} \cup \mathcal{W}^-$$

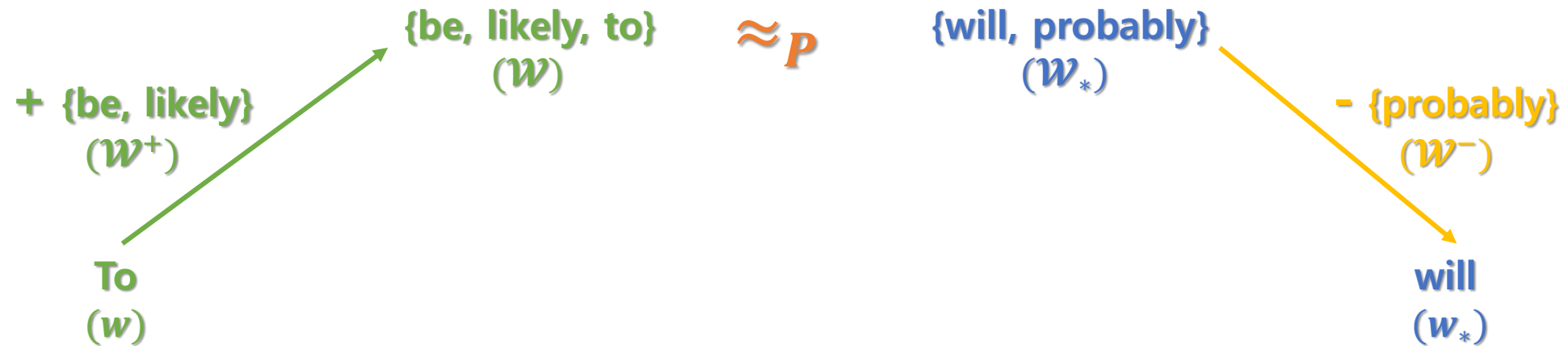


**Proof**

-Analogy

**<Word Transformation>**

$$\mathcal{W} = \{w_x\} \cup \mathcal{W}^+, \mathcal{W}_* = \{w_{x_*}\} \cup \mathcal{W}^-$$

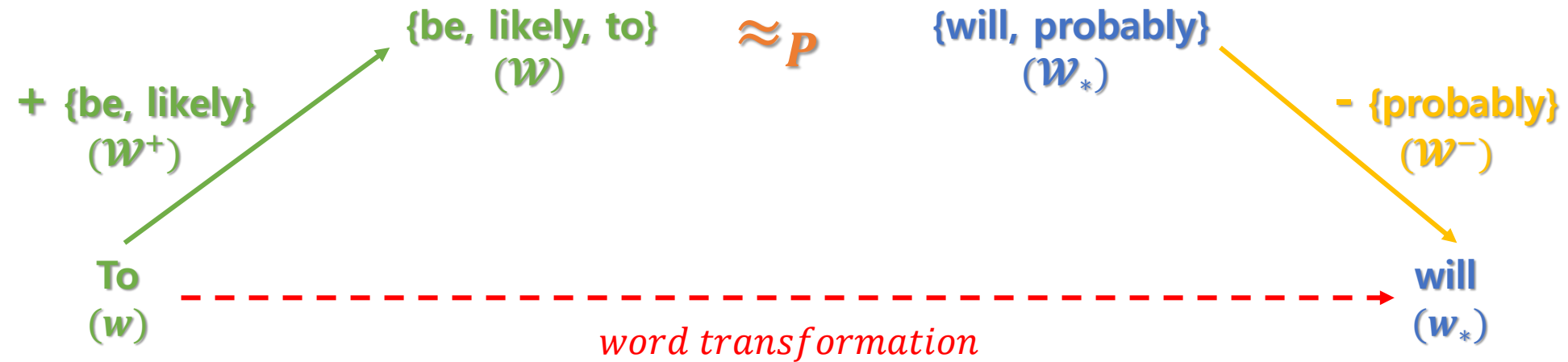


**Proof**

-Analogy

**<Word Transformation>**

$$\mathcal{W} = \{w_x\} \cup \mathcal{W}^+, \mathcal{W}_* = \{w_{x_*}\} \cup \mathcal{W}^-$$

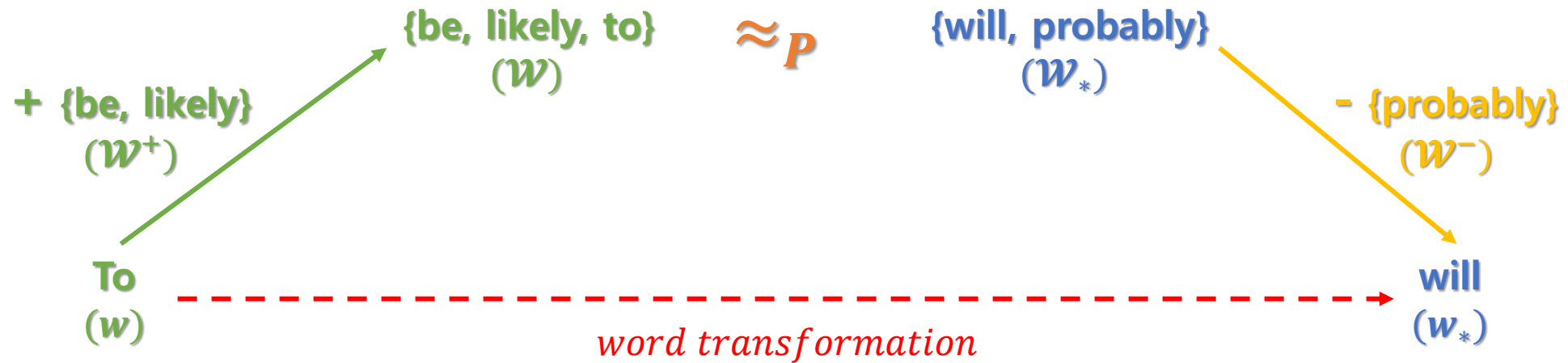


**Proof**

-Analogy

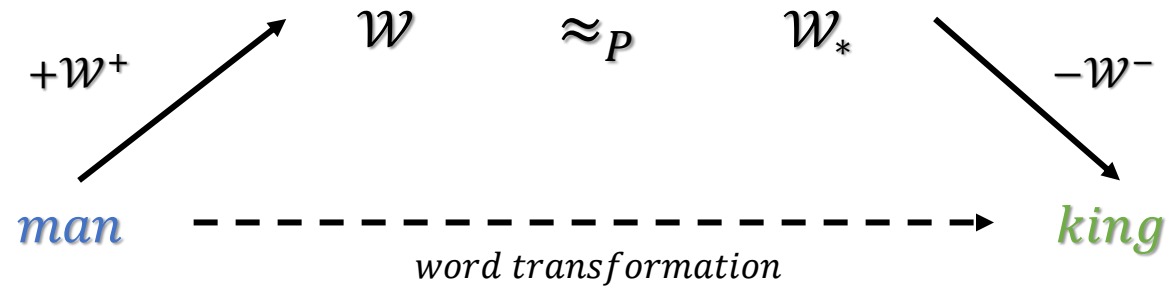
## &lt;Word Transformation&gt;

**Definition D3.** There exists a **word transformation** from  $w_x \in \mathcal{E}$  to  $w_{x^*} \in \mathcal{E}$  with **transformation parameters**  $\mathcal{W}^+, \mathcal{W}^- \subseteq \mathcal{E}$  iff  $\{w_x\} \cup \mathcal{W}^+ \approx_P \{w_{x^*}\} \cup \mathcal{W}^-$



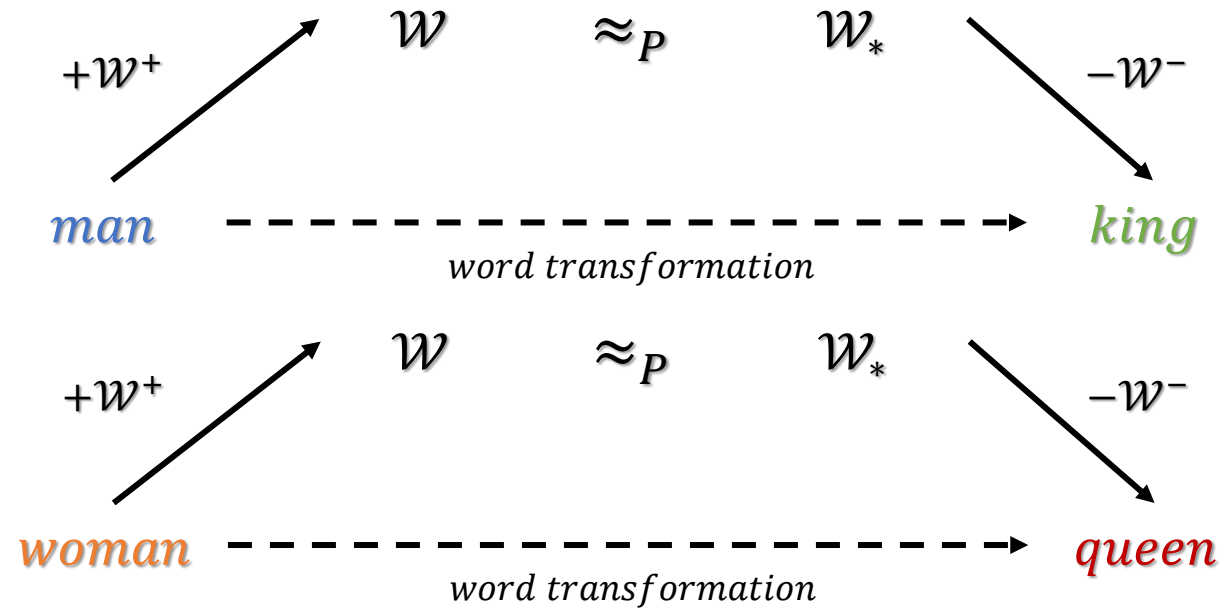
**Proof**

-Analogy

**<Word Transformation to Analogy>**

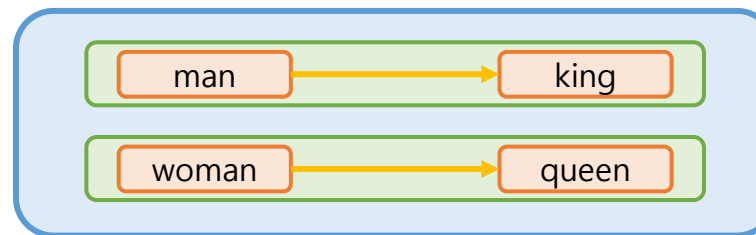
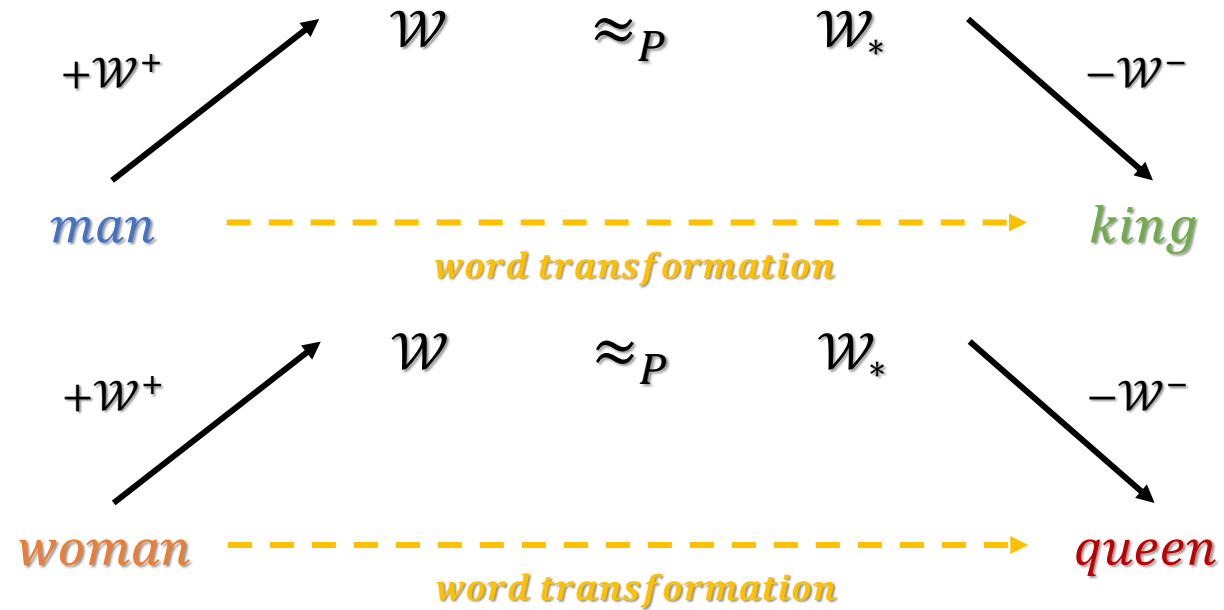
**Proof**

-Analogy

**<Word Transformation to Analogy>**

**Proof**

-Analogy

**<Word Transformation to Analogy>**

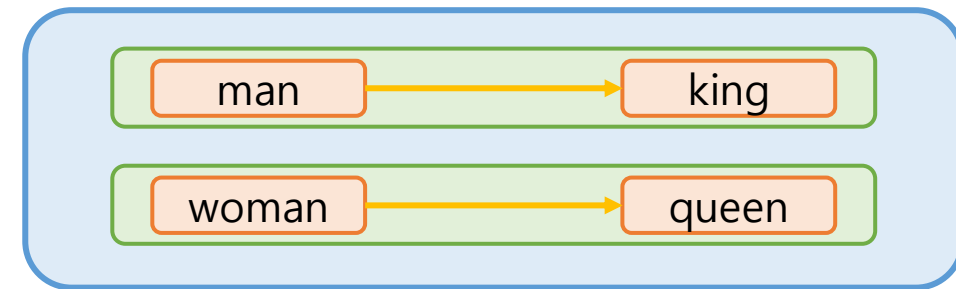
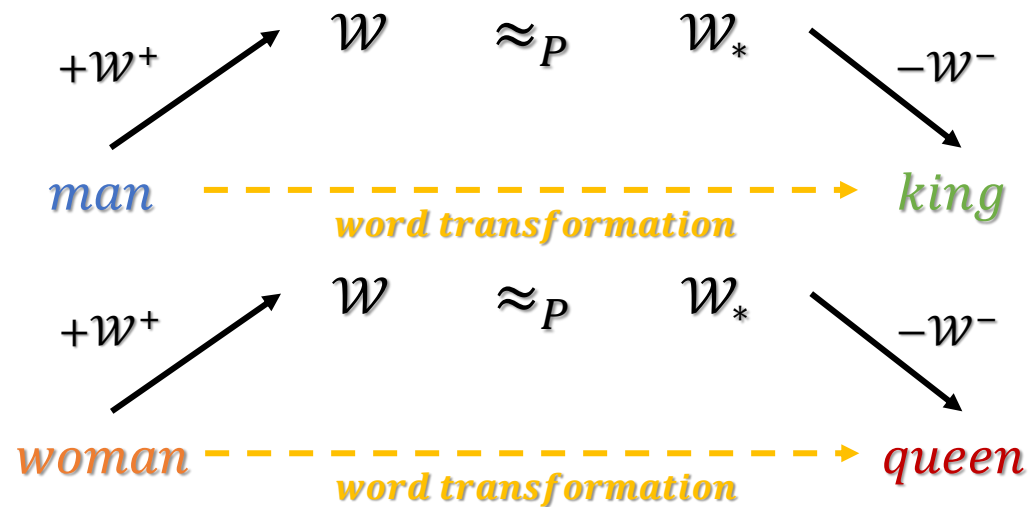


**Proof**

-Analogy

**<Word Transformation to Analogy>**

**Definition D4.** We say “ $w_a$  is to  $w_{a^*}$  as  $w_b$  is to  $w_{b^*}$ ” for  $w_a, w_b, w_{a^*}, w_{b^*} \in \mathcal{E}$  iff there exist parameters  $\mathcal{W}^+, \mathcal{W}^- \subseteq \mathcal{E}$  that **simultaneously transform**  $w_a$  to  $w_{a^*}$  and  $w_b$  to  $w_{b^*}$



**Proof**

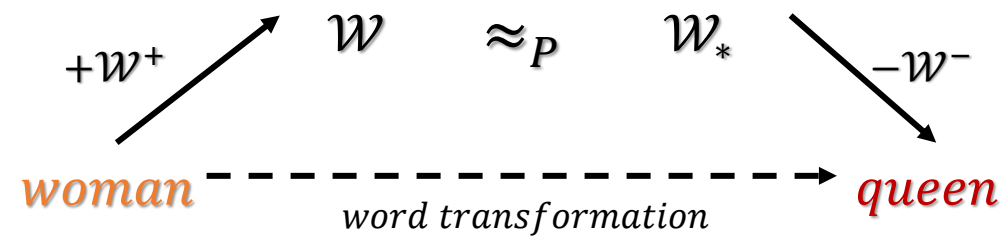
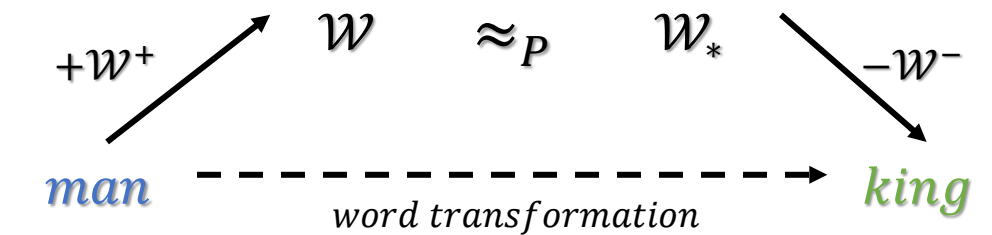
-Analogy

**<Analogy to Paraphrase>**

That is, we say:

"**man** is to **king** as **woman** is to **queen**"

*iff* there exist parameters  $\mathcal{W}^+, \mathcal{W}^- \subseteq \mathcal{E}$  that simultaneously transform **man** to **king** and **woman** to **queen**



**Proof**

-Analogy

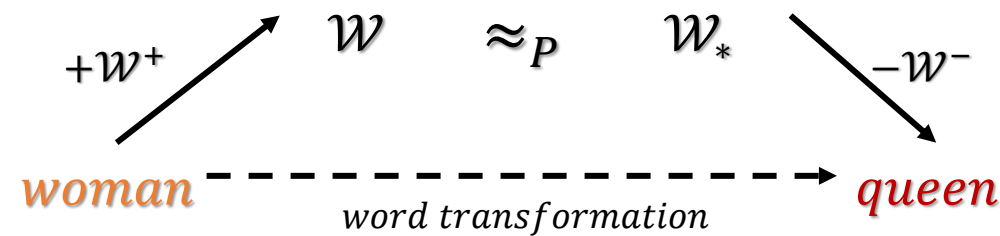
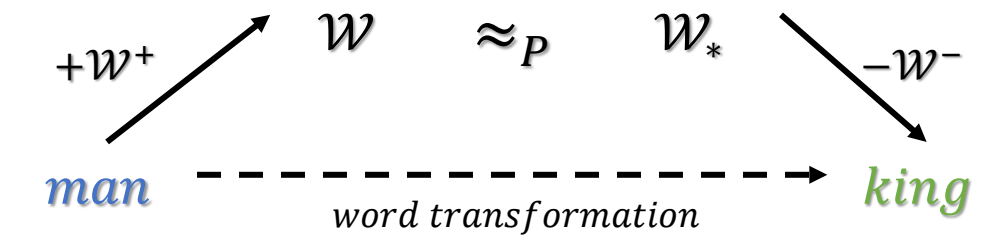
**<Analogy to Paraphrase>**

That is, we say:

"**man** is to **king** as **woman** is to **queen**"

*iff* there exist parameters  $\mathcal{W}^+, \mathcal{W}^- \subseteq \mathcal{E}$  that simultaneously transform **man** to **king** and **woman** to **queen**

Let  $\mathcal{W}^+ = \{\text{king}\}$ ,  
 $\mathcal{W}^- = \{\text{man}\}$



**Proof**

-Analogy

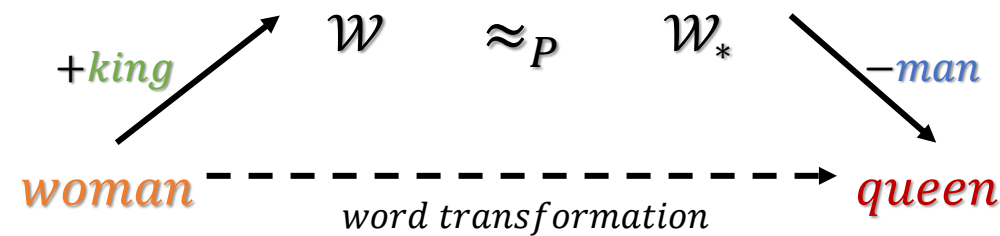
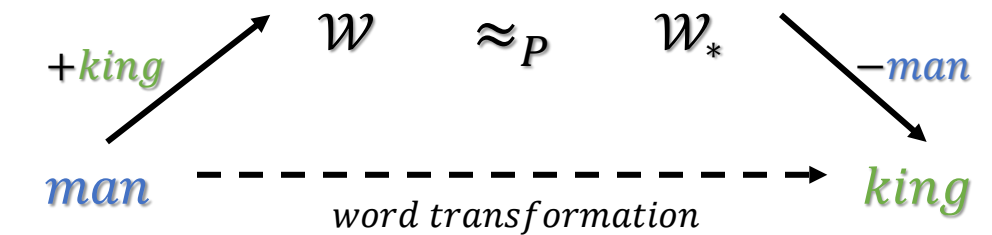
**<Analogy to Paraphrase>**

That is, we say:

"**man** is to **king** as **woman** is to **queen**"

*iff* there exist parameters  $\mathcal{W}^+, \mathcal{W}^- \subseteq \mathcal{E}$  that simultaneously transform **man** to **king** and **woman** to **queen**

Let  $\mathcal{W}^+ = \{\text{king}\}$ ,  
 $\mathcal{W}^- = \{\text{man}\}$



**Proof**

-Routemap

**<Routemap>**

"man is to king as woman is to queen"



man *transforms to* king as woman *transforms to* queen



{woman, king} *paraphrases* {man, queen}



*Dependency  
Error*

$$\text{PMI}_{\text{king}} - \text{PMI}_{\text{man}} + \text{PMI}_{\text{woman}} \approx \text{PMI}_{\text{queen}}$$



$$\text{PMI}_i \approx \mathbf{w}_i^T \mathbf{C}$$

$$\mathbf{w}_{\text{king}} - \mathbf{w}_{\text{man}} + \mathbf{w}_{\text{woman}} \approx \mathbf{w}_{\text{queen}}$$

semantic

geometric

## <Analogies Explained>

**Lemma 3.** If “ $w_a$  is to  $w_{a^*}$  as  $w_b$  is to  $w_{b^*}$ ” by **D4** with transformation parameters  $\mathcal{W}^+, \mathcal{W}^- \subseteq \mathcal{E}$ , then:

$$\begin{aligned}
 PMI_{b^*} &= PMI_{a^*} - PMI_a + PMI_b \\
 &+ \rho^{\mathcal{W}^b, \mathcal{W}_*^b} - \rho^{\mathcal{W}^a, \mathcal{W}_*^a} \\
 &+ \left( \sigma^{\mathcal{W}^b} - \sigma^{\mathcal{W}_*^b} \right) - \left( \sigma^{\mathcal{W}^a} - \sigma^{\mathcal{W}_*^a} \right) \\
 &- \left( \left( \tau^{\mathcal{W}^b} - \tau^{\mathcal{W}_*^b} \right) - \left( \tau^{\mathcal{W}^a} - \tau^{\mathcal{W}_*^a} \right) \right) 1
 \end{aligned}$$

where  $\mathcal{W}^x = \{w_x\} \cup \mathcal{W}^+$ ,  $\mathcal{W}_*^x = \{w_{x^*}\} \cup \mathcal{W}^-$  for  $x \in \{a, b\}$  and  $\rho^{\mathcal{W}^b, \mathcal{W}_*^b}, \rho^{\mathcal{W}^a, \mathcal{W}_*^a}$  are small.

**Proof of Lemma 3.** Let  $\mathcal{W} = \mathcal{W}^x$ ,  $\mathcal{W}_* = \mathcal{W}_*^x$  for  $x \in \{a, b\}$  in instance of **Cor 2.1** and take the difference.  $\mathcal{W}^x$  paraphrases  $\mathcal{W}_*^x$  for  $x \in \{a, b\}$  by **D3** and **D4**

**Proof**

-Overall

**<Analogies Explained>**

**Theorem 3.** (Analogies) If “ $w_a$  is to  $w_{a^*}$  as  $w_b$  is to  $w_{b^*}$ ” by **D4** with transformation parameters  $\mathcal{W}^+, \mathcal{W}^- \subseteq \mathcal{E}$ , then:

$$\begin{aligned} w_{b^*} = & w_{a^*} - w_a + w_b \\ & + C^\dagger (\rho^{\mathcal{W}^b, \mathcal{W}_*^b} - \rho^{\mathcal{W}^a, \mathcal{W}_*^a} \\ & + (\sigma^{\mathcal{W}^b} - \sigma^{\mathcal{W}_*^b}) - (\sigma^{\mathcal{W}^a} - \sigma^{\mathcal{W}_*^a}) \\ & - ((\tau^{\mathcal{W}^b} - \tau^{\mathcal{W}_*^b}) - (\tau^{\mathcal{W}^a} - \tau^{\mathcal{W}_*^a}))1) \end{aligned}$$

“man is to king as woman is to queen”

$\Downarrow (\rho, \sigma, \tau)$

$$w_{king} - w_{man} + w_{woman} \approx w_{queen}$$

**Proof**

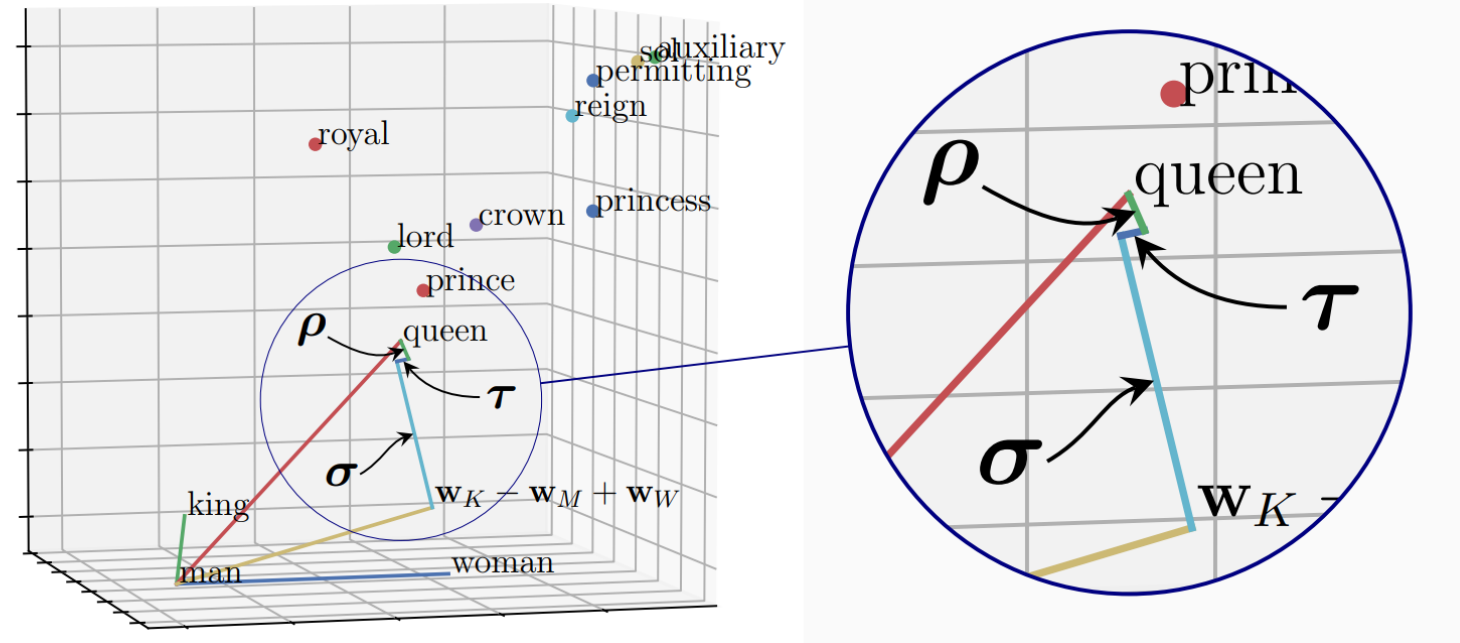
-Overall

## &lt;Analogies Explained&gt;

"man is to king as woman is to queen"

$\Downarrow (\rho, \sigma, \tau)$

$$W_{\text{king}} - W_{\text{man}} + W_{\text{woman}} \approx W_{\text{queen}}$$





# Conclusion

### <Conclusion>

- To derive a probabilistic definition of ***paraphrasing*** and show that it governs the relationship between one word embedding and any sum of others
- To show how paraphrasing can be generalized and interpreted as the ***transformation*** from one word to another, giving a mathematical formulation for “ $w_x$  is to  $w_{x^*}$ ”
- To provide the first rigorous proof of the linear relationship between word embeddings of analogies, including explicit, interpretable error terms

**Any Questions?**

**Thank You**