

word representations umumnya vector

(Distributed Word Representations) Word Embeddings

Alfan F. Wicaksono
Information Retrieval Lab.
Faculty of Computer Science
Universitas Indonesia

How to more robustly match a user's intent?

Query: “motel di depok”

Sistem diharapkan tidak hanya mengembalikan dokumen tentang “motel di depok”; tetapi juga “hotel di depok”, “penginapan di depok”, dsb.

Expansion Using Query Logs

Context-Free Query Expansion

- Misal, dari analisis query logs didapatkan bahwa **“wet ground”** = **“wet earth”**.
- Oleh karena itu, jika kedepannya ada query term **“ground”**, perlu di-expand dengan term **“earth”**.
- Problem: **“ground coffee”** -> **“earth coffee”** ??

Thesaurus-based Approach

setiap kata ada hierarki

Contoh: Path-based similarity

Two terms are similar if they are near each other in the thesaurus hierarchy.

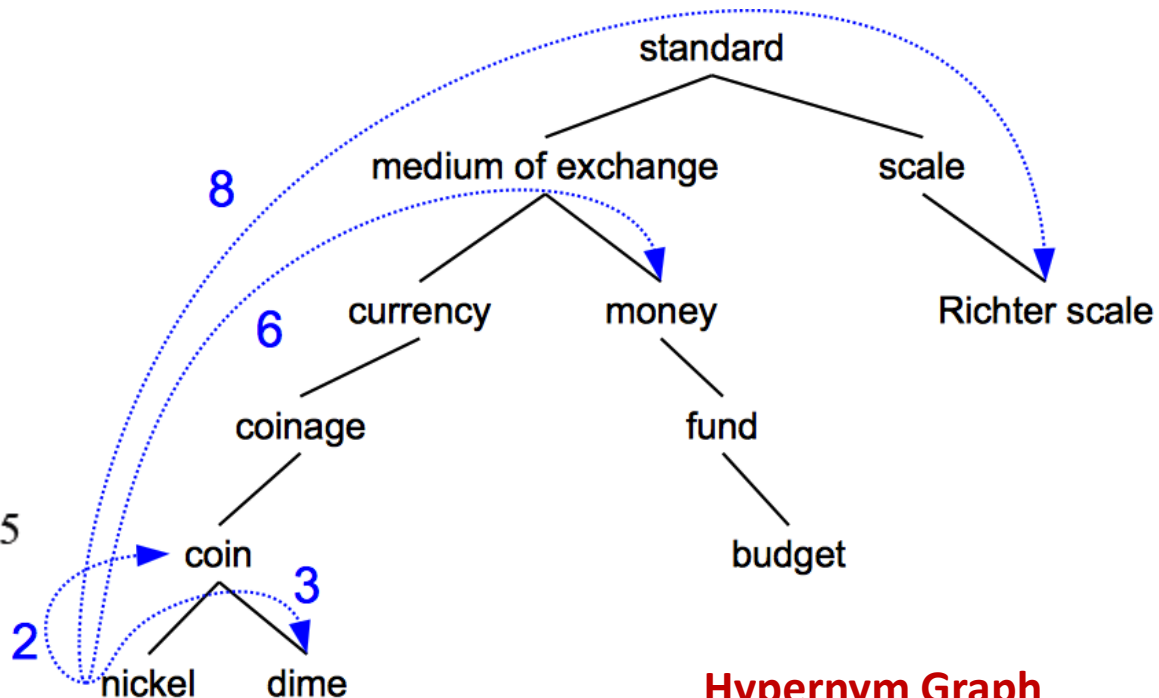
$$\text{simpath}(c_1, c_2) = \frac{1}{\text{pathlen}(c_1, c_2)}$$

$$\text{wordsim}(w_1, w_2) = \max_{c_1 \in \text{senses}(w_1), c_2 \in \text{senses}(w_2)} \text{sim}(c_1, c_2)$$

$$\text{simpath}(\text{nickel}, \text{coin}) = 1/2 = .5$$

$$\text{simpath}(\text{fund}, \text{budget}) = 1/2 = .5$$

$$\text{simpath}(\text{nickel}, \text{currency}) = 1/4 = .25$$



Hypernym Graph

Thesaurus-based Approach

Problems...

- We don't have a thesaurus for every language
- For Indonesian, our WordNet is not complete
 - Many words are missing
 - Connections between senses are missing
 - ...

Word -> Vector; Term Relations?

Pada konsep VSM standar, setiap term menjadi **basis/axis** di vector space.

Query: “hotel”

kalau pakai tf idf aja gak mungkin cocok

Document: “motel motel motel”

$$q = [0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, \dots, 0]$$

$$d = [0, 0, 3, 0, 0, 0, 0, 0, 0, 0, 0, \dots, 0]$$

$$\text{sim}(q, d) = q \cdot d^T = 0$$

Orthogonal!

Word -> Vector; Term Relations?

Pada konsep VSM standar, setiap term menjadi **basis/axis** di vector space.

$$\text{sim}(q, d) = q \cdot d^T$$

q = "hotel di depok"

$$\begin{matrix} & [1, 0, 1, 1] & \times & \begin{bmatrix} 0 \\ 1 \\ 1 \\ 1 \end{bmatrix} & = & 2.0 & \text{Make Sense?} \\ \swarrow & \uparrow & \uparrow & \nwarrow & & & \\ \text{hotel} & \text{motel} & \text{di} & \text{depok} & & & \end{matrix}$$

d = "motel di depok"

Orthogonal!

Word -> Vector; Term Relations?

Query Translation Model (Berger & Lafferty, 99)

Tambahkan **trainable parameter W** untuk menangkap **similarity** atau **translasi** antar kata.

$$W = V * V$$

ini adalah W (trainable parameter)

	hotel	motel	di	depok
hotel	1.0	0.8	0.0	0.0
motel	0.8	1.0	0.0	0.0
di	0.0	0.0	1.0	0.0
depok	0.0	0.0	0.0	1.0

Word -> Vector; Term Relations?

Query Translation Model (Berger & Lafferty, 99)

Tambahkan parameter **W** untuk menangkap similarity antar kata.

$$\text{sim}(q, d) = q \cdot W \cdot d^T$$

q = "hotel di depok"

$$\begin{bmatrix} 1 & 0 & 1 & 1 \end{bmatrix} \times \begin{bmatrix} 1.0 & 0.8 & 0.0 & 0.0 \\ 0.8 & 1.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 1.0 & 0.0 \\ 0.0 & 0.0 & 0.0 & 1.0 \end{bmatrix} \times \begin{bmatrix} 0 \\ 1 \\ 1 \\ 1 \end{bmatrix} = 2.8$$

Better!

d = "motel di depok"

Tetap saja isu Sparsity muncul. Ukuran W besar!

Distributional-based Approach

- Can we learn a **dense low-dimensional** representation of a word such that dot products $u \cdot v^T$ express word similarity?
- Masih mungkin juga menyisipkan “translation” matrix antar vocab (misal, cross-language): $u \cdot W \cdot v^T$

Distributional-based Approach

- Based on the idea that contextual information alone constitutes a viable representation of linguistic items.
- As opposed to formal linguistics and the Chomsky tradition.

lihat context based on surrounding terms

- Zellig Haris (1954): “...if A and B have almost identical environments, we say that they are synonyms...”
- J. R. Firth (1957): “You shall know a word by the company it keeps”

Distributional-based Approach

gogos

Ada yang tahu apa itu **gogos**?

Distributional-based A



Makan **gogos** dengan sambal sungguh nikmat.
Beras ketan diperlukan untuk membuat **gogos**.
Teman-teman menikmati **gogos** hangat di kantin.
Menikmati makanan **gogos**, lemper dari makasar.

- From context words humans can guess **gogos** means
 - A traditional food
- Intuition for algorithm:
 - **Two words are similar if they have similar word contexts**

Latent Semantic Analysis

Alin Revisited: Rank

Let C be an $M \times N$ matrix. The **rank** of a matrix is the number of **linearly independent rows** (or columns) in it; thus, $\text{rank}(C) \leq \min\{M, N\}$.

$$C = \begin{bmatrix} 2 & 4 & 8 \\ 1 & 2 & 4 \end{bmatrix}, \quad \text{rank}(C) = 1$$

$$C = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}, \quad \text{rank}(C) = 3$$

full rank

Alin Revisited: Rank

Let C be an $M \times N$ matrix. The **rank** of a matrix is the number of **linearly independent rows** (or columns) in it; thus, $\text{rank}(C) \leq \min\{M, N\}$.

Secara umum, ubah dahulu ke bentuk **Row-Echelon Form**; lalu hitung ada **berapa baris yang non-zero**.

$$C = \begin{bmatrix} 1 & 2 & 3 \\ 2 & 1 & 4 \\ 3 & 0 & 5 \end{bmatrix} \xrightarrow{\quad} \begin{bmatrix} 1 & 2 & 3 \\ 0 & -3 & -2 \\ 0 & -6 & -4 \end{bmatrix} \xrightarrow{\quad} \begin{bmatrix} 1 & 2 & 3 \\ 0 & -3 & -2 \\ 0 & 0 & 0 \end{bmatrix}$$

$$\begin{aligned} R_2 &\rightarrow R_2 - 2R_1 \\ R_3 &\rightarrow R_3 - 3R_1 \end{aligned}$$

$$R_3 \rightarrow R_3 - 2R_2$$

$$\text{Rank}(C) = 2$$

Alin Revisited: Eigenvector

For a square $M \times M$ matrix C and a non-zero vector \mathbf{x} , the values of λ satisfying

$$C\mathbf{x} = \overset{\text{eigen value}}{\lambda} \underset{\text{eigen vector}}{\mathbf{x}}$$

are called the **eigenvalues** of C .

\mathbf{x} disebut **eigenvector**.

Banyaknya non-zero eigenvalues dari C adalah paling banyak $\text{Rank}(C)$.

Contoh:

$$C = \begin{bmatrix} -6 & 3 \\ 4 & 5 \end{bmatrix}$$

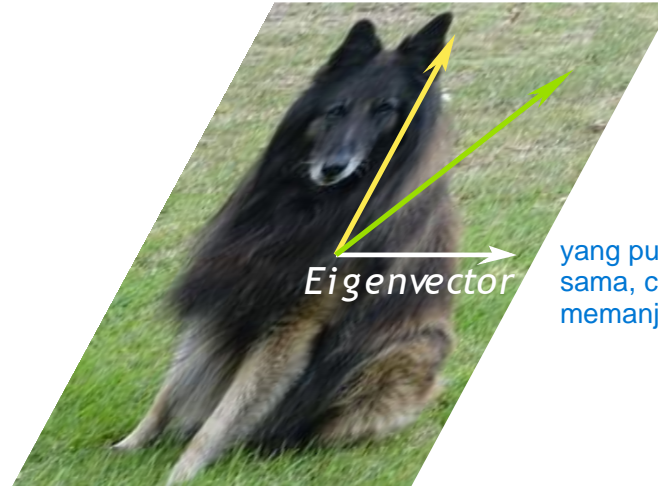
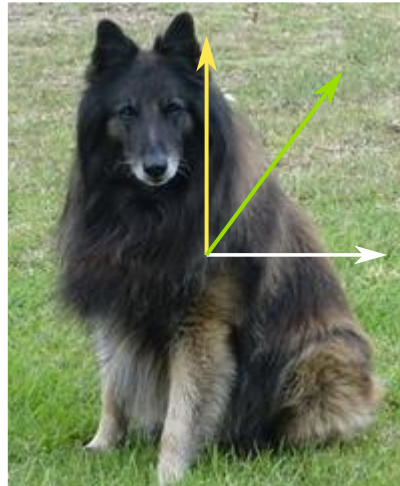
$$\mathbf{x}_1 = \begin{bmatrix} 1 \\ 4 \end{bmatrix} \quad \lambda_1 = 6$$

$$\mathbf{x}_2 = \begin{bmatrix} -3 \\ 1 \end{bmatrix} \quad \lambda_2 = -7$$

Alin Revisited: Eigenvector

Singkat cerita, eigenvector tidak akan berubah arah ketika ditransformasi (hanya memanjang atau memendek).

“ditransformasi” = “dikali matriks”



yang putih eigen vector arahnya sama, cuman bisa memanjang/memendek.

Alin Revisited: Eigenvector

Sekedar istilah ...

Right-eigenvector of C

$$Cx = \lambda x$$

$$C = \begin{bmatrix} -6 & 3 \\ 4 & 5 \end{bmatrix} \quad x = \begin{bmatrix} 1 \\ 4 \end{bmatrix} \quad \lambda = 6$$

Left-eigenvector of C

$$y^T C = \lambda y^T$$

$$C = \begin{bmatrix} -6 & 3 \\ 4 & 5 \end{bmatrix} \quad y^T = [\ ? \quad ?] \quad \lambda = ?$$

Latihan: cari sebuah left-eigenvector dari C!

Vector Space Model

Term-Document Matrix

Each cell is the count of word t in document d (bisa juga TF-IDF)

	D1	D2	D3	D4	D5
ekonomi	0	1	40	38	1
pusing	4	5	1	3	30
keuangan	1	2	30	25	2
sakit	4	6	0	4	25
inflasi	8	1	15	14	1

document embedding

Vector of D2 = [1, 5, 2, 6, 1]

Vector Space Model

Term-Document Matrix

Each cell is the count of word t in document d (bisa juga TF-IDF)

	D1	D2	D3		D4	D5
ekonomi	0	1	40		38	1
pusing	4	5	1		3	30
keuangan	1	2	30		25	2
sakit	4	6	0		4	25
inflasi	8	1	15		14	1

Two documents are similar if they have similar vector!

D3 = [40, 1, 30, 0, 15]

D4 = [38, 3, 25, 4, 14]

Vector Space Model

Term-Document Matrix

Each cell is the count of word t in document d (bisa juga TF-IDF)

	D1	D2	D3	D4	D5
ekonomi	0	1	40	38	1
pusing	4	5	1	3	30
keuangan	1	2	30	25	2
sakit	4	6	0	4	25
inflasi	8	1	15	14	1

word
embeddin
g

Vector of word “sakit” = [4, 6, 0, 4, 25]

word embedding: document yang jadi context

However, kalau document embedding: word-wordnya yang jadi context

Vector Space Model

Term-Document Matrix

Each cell is the count of word t in document d (bisa juga TF-IDF)

	D1	D2	D3	D4	D5
ekonomi	0	1	40	38	1
pusing	4	5	1	3	30
keuangan	1	2	30	25	2
sakit	4	6	0	4	25
inflasi	8	1	15	14	1

Two words are similar if they have similar vector!

pusing = [4, 5, 1, 3, 30]

sakit = [4, 6, 0, 4, 25]

Vector Space Model

Jika **C** adalah term-document matrix,

contoh C = term ada 3, document ada 4 jadi 3x4. Nah hasilnya bakal 3x3 dimana setiap elemen itu menandakan kemiripan antar term (term_1 dengan term_2)

$C.C^T$ mengandung dot products (similarity) dari semua pasangan **term vectors**; dan

- 1. matrixsnya persegi
- 2. matriksnya simetrik

kebalikannya, ini menandakan kemiripan antar document

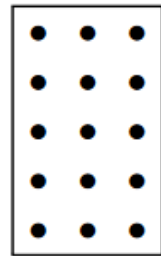
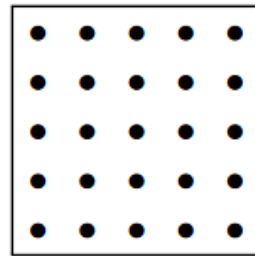
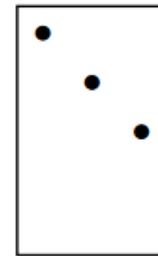
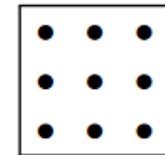
$C^T.C$ mengandung dot products (similarity) dari semua pasangan **document vectors**.

Dekomposisi Term-Document Matrix!

Misal, C adalah term-document matrix. C didekomposisi menjadi perkalian 3 matrix (**Singular Value Decomposition**).

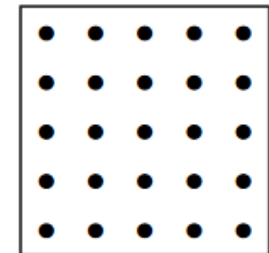
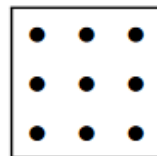
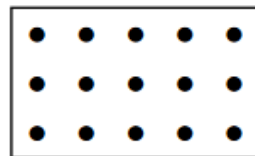
$$C = U \times \Sigma \times V^T \quad C \in R^{M \times N}$$

Jika $M > N$

 C $=$  U  Σ  V^T

Cell selain titik = 0

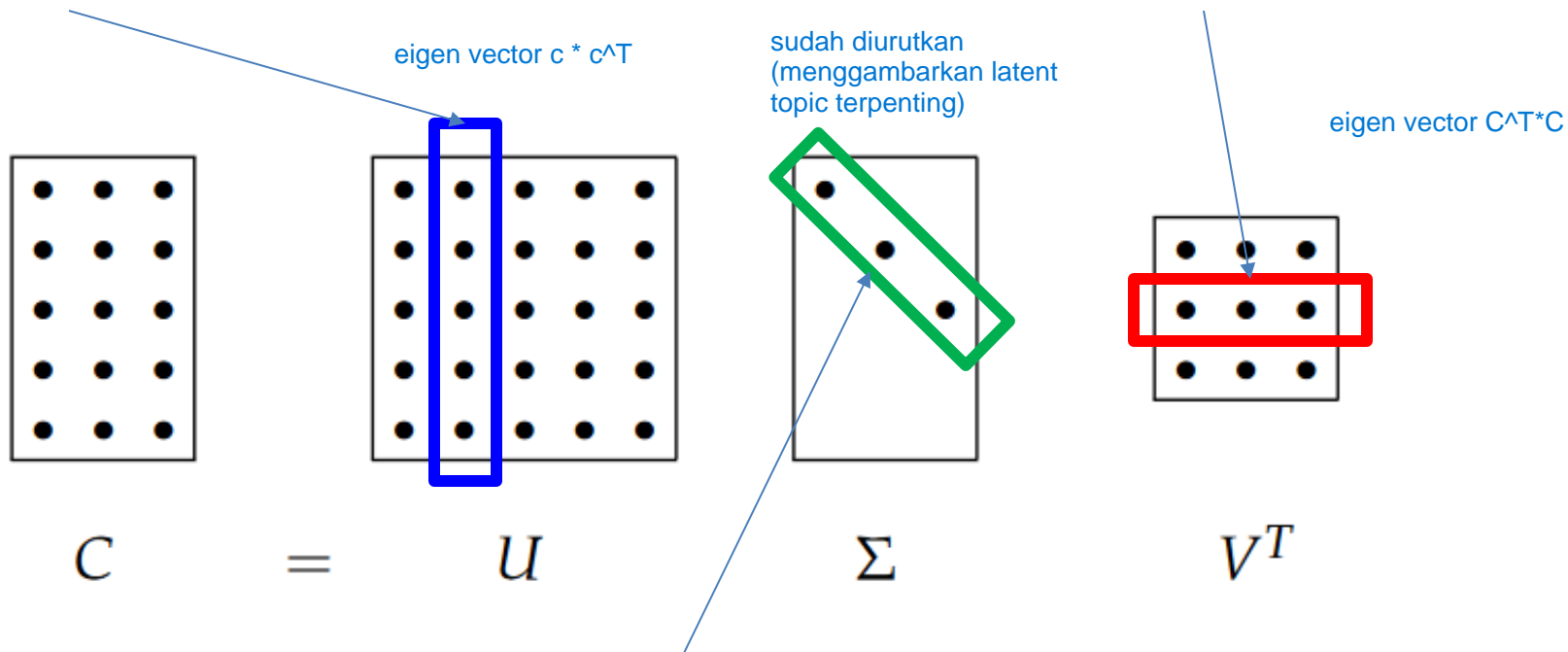
Jika $M < N$



Cara menghitung U , Σ , dan V^T

Kolom-kolom pada U adalah **orthogonal eigenvector** dari $C.C^T$

Baris-baris pada V^T adalah **orthogonal eigenvector** dari $C^T.C$



Singular Values (**Akar kuadrat dari Eigenvalues**) dari $C^T.C$ atau $C.C^T$ (sama saja).

Interpretasi Truncated SVD

Misal, C adalah term-document matrix. C didekomposisi menjadi perkalian 3 matrix.

$$C = U \times \Sigma \times V^T$$

ortogonal

Kolom yang di U , maksudnya topic 1, adalah eigen vector dari $C \times C^T$

barisnya itu term
kolom itu document (menjadi context)

Term Document Matrix

	Doc-1	Doc-2	Doc-3	Doc-4
Term-1				
Term-2				
Term-3				
Term-4				

baris term
kolom: latent topic yang jadi context bagi term

Word Assignment to Topics

	Topic-1	Topic-2
Term-1		
Term-2		
Term-3		
Term-4		

seberapa penting latent topic

Topic Importance

	Topic-1	Topic-2
Topic-1		
Topic-2		

SIGMA (matriks diagonal)

V^T :
- kolom: embedding document terbaru
Misal document embedding dokumen 1 terhadap topic 1

Topic Distribution Across Documents

	Doc-1	Doc-2	Doc-3	Doc-4
Topic-1				
Topic-2				

ini matriks ortogonal juga

topic: campuran dokumen-dokumen

U (Matriks ortogonal)

kolomnya itu $C^T \times C$

kita melakukan diagonalisasi terhadap matriks covariant untuk menghilangkan regularisasi (situasi dimana fitur X korelasinya tinggi dengan fitur Y)

Gambar: <https://www.datacamp.com/tutorial/discovering-hidden-topics-python>

	D1	D2	D3	D4	D5
Cancer	6	0	10	1	7
Flower	2	8	1	9	0
Tumor	6	2	7	0	8
Rose	1	6	0	7	1

word embedding kata cancer yang lama

C

U ini ortogonal kalau dihitung L2 Norm hasilnya 1 (artinya setiap term itu tegak lurus)

diurutin sesuai kepentingan. Lihat deh besarnya signifikan.
Artinya ada 2 latent topic yang signifikan

Σ	18.93	0	0	0
	0	14.49	0	0
	0	0	2.60	0
	0	0	0	0.86

	1	2	3	4
Cancer	0.66	0.33	0.64	0.18
Flower	0.33	-0.71	0.18	-0.57
Tumor	0.61	0.25	-0.72	-0.19
Rose	0.24	-0.55	-0.18	0.77

	D1	D2	D3	D4	D5
1	0.45	0.28	0.59	0.28	0.51
2	0.10	-0.59	0.30	-0.69	0.26
3	-0.11	-0.41	0.59	0.38	-0.56
4	-0.51	-0.42	-0.12	0.44	0.58
5	-0.70	0.46	0.42	-0.30	0.04

$U * S$: word embedding kata yang baru

$S * V^T$ = document embedding yang baru

U

di V^T setiap dokumen-dokumennya ini ortogonal

Kalau dihitung L2 Norm itu nilainya 1

V^T

```
import numpy as np; u, s, vt = np.linalg.svd(C, full_matrices = True)
```

LSA: Low-rank approximation of C

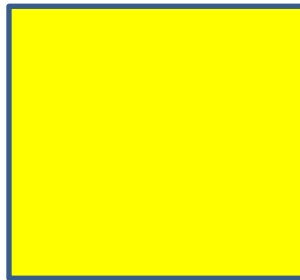
Truncating U , Σ , V^T to **K dimensions** produces best possible **K** rank approximation of original matrix C .

Buang (jadikan nol) **L - K singular values paling kecil** di Σ !

$$U, M \times L$$



$$\Sigma, L \times L$$



$$V^T, L \times N$$

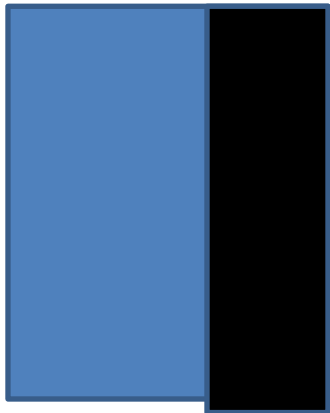


C_k = rank-k approximation of C

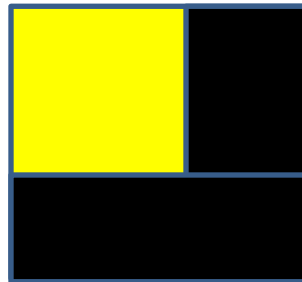
Truncating U , Σ , V^T to **K dimensions** produces best possible **K** rank approximation of original matrix **C**.

Jadikan NOL **L - K** singular values paling kecil di Σ !

$$U_k, M \times K$$



$$\Sigma_k, K \times K$$



$$V_k^T, K \times N$$



$$C_k = U_k \times \Sigma_k \times V_k^T$$

	D1	D2	D3	D4	D5
Cancer	6	0	10	1	7
Flower	2	8	1	9	0
Tumor	6	2	7	0	8
Rose	1	6	0	7	1

C

Σ

18.93	0	0	0
0	14.49	0	0
0	0	2.60	0
0	0	0	0.86

	1	2	3	4
Cancer	0.66	0.33	0.64	0.18
Flower	0.33	-0.71	0.18	-0.57
Tumor	0.61	0.25	-0.72	-0.19
Rose	0.24	-0.55	-0.18	0.77

U

	D1	D2	D3	D4	D5
1	0.45	0.28	0.59	0.28	0.51
2	0.10	-0.59	0.30	-0.69	0.26
3	-0.11	-0.41	0.59	0.38	-0.56
4	-0.51	-0.42	-0.12	0.44	0.58
5	-0.70	0.46	0.42	-0.30	0.04

V^T

```
import numpy as np; u, s, vt = np.linalg.svd(C, full_matrices = True)
```

Contoh: rank-2
approximation of C

	D1	D2	D3	D4	D5
Cancer	6.27	0.76	9.03	0.29	7.85
Flower	1.80	7.99	0.65	9.04	0.57
Tumor	5.70	1.15	8.09	0.79	7.03
Rose	1.29	6.08	0.38	6.89	0.33

C_2

Σ_2

18.93	0
0	14.49

	1	2
Cancer	0.66	0.33
Flower	0.33	-0.71
Tumor	0.61	0.25
Rose	0.24	-0.55

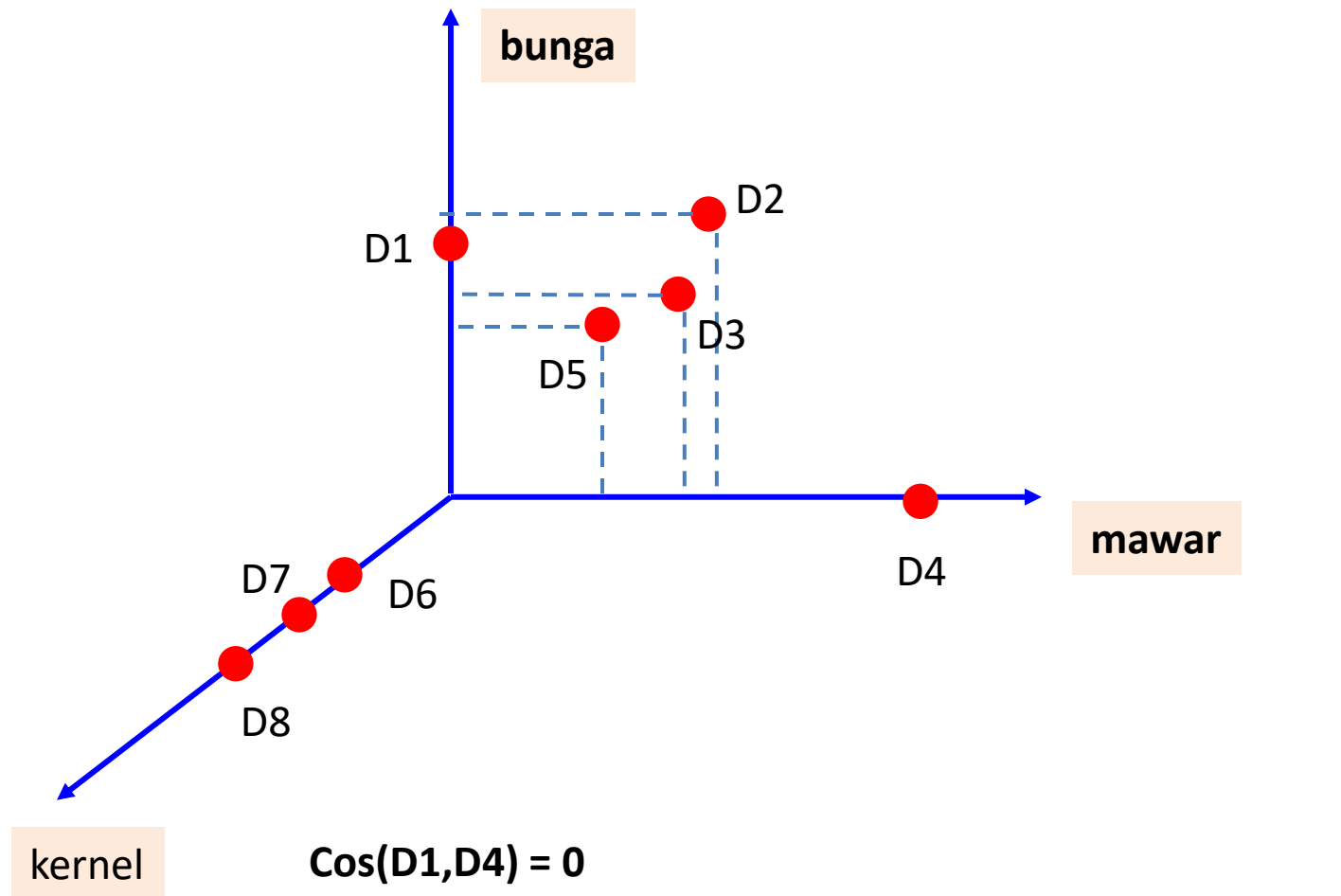
U_2

	D1	D2	D3	D4	D5
1	0.45	0.28	0.59	0.28	0.51
2	0.10	-0.59	0.30	-0.69	0.26

V_2^T

```
import numpy as np; u, s, vt = np.linalg.svd(C, full_matrices = True)
```


Jadi mengapa LSA/SVD berhasil menangkap “similarity”
antar kata?

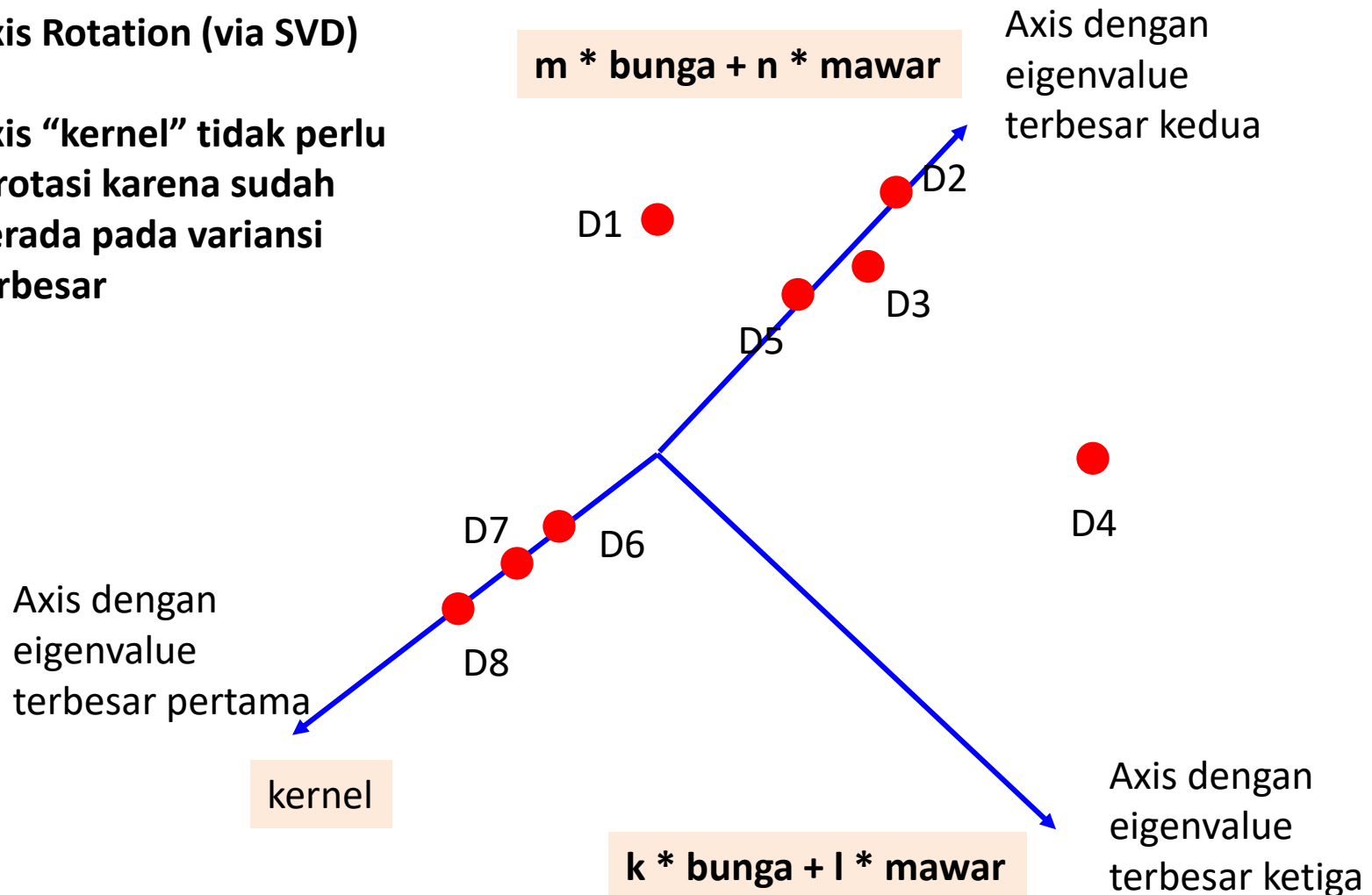


Padahal “bunga” dan “mawar” secara semantic “dekat”

Jadi mengapa LSA/SVD berhasil menangkap “similarity” antar kata?

Axis Rotation (via SVD)

Axis “kernel” tidak perlu dirotasi karena sudah berada pada variansi terbesar



Jadi mengapa LSA/SVD berhasil menangkap “similarity” antar kata?

Buang Axis dengan eigenvalue paling kecil.

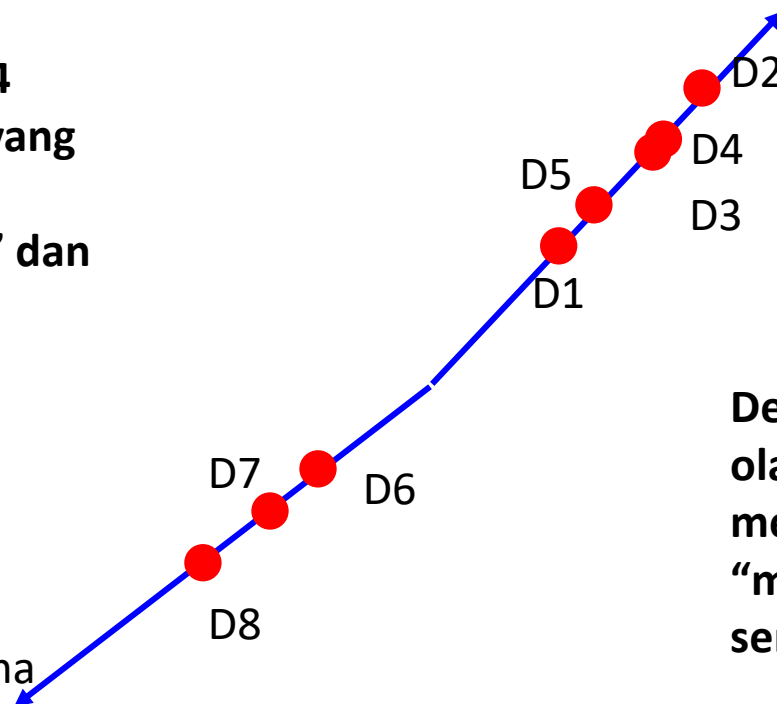
Akibatnya D1 & D4 berada pada axis yang sama, yaitu axis gabungan “bunga” dan “mawar”.

Axis dengan eigenvalue terbesar pertama

kernel

$m * \text{bunga} + n * \text{mawar}$

Axis dengan eigenvalue terbesar kedua



Dengan begitu, seolah-olah LSA/SVD berhasil menemukan “bunga” dan “mawar” dekat secara semantik.

$$\text{Cos}(D1, D4) > 0$$

Dimanakah *Document Vector*?

Jika ingin mendapatkan **rank-k vector representation** dari sebuah dokumen di **original vector space**:

=> Vektor kolom pada C_k

Dari Contoh sebelumnya:

	D1	D2	D3	D4	D5
Cancer	6.27	0.76	9.03	0.29	7.85
Flower	1.80	7.99	0.65	9.04	0.57
Tumor	5.70	1.15	8.09	0.79	7.03
Rose	1.29	6.08	0.38	6.89	0.33

C_2

Rank-2 Document Embedding dari D3 = [9.03, 0.65, 8.09, 0.38]

Dimanakah *Document Vector*?

Jika ingin mendapatkan **rank-k vector representation** dari sebuah dokumen di **vector space baru (rank-k subspace)**:

=> Vektor kolom pada $\Sigma_k \times V_k^T$

Dari Contoh sebelumnya:

2-Dimensional Document Embedding dari D3 adalah

$$\Sigma_2 \begin{array}{|c|c|} \hline 18.93 & 0 \\ \hline 0 & 14.49 \\ \hline \end{array} \times V_2^T \begin{array}{|c|} \hline \text{D3} \\ \hline 0.59 \\ \hline 0.30 \\ \hline \end{array} = \begin{array}{|c|} \hline \text{D3} \\ \hline 11.17 \\ \hline 4.35 \\ \hline \end{array}$$

Dimanakah *Term Vector (Word Embedding)*?

Jika ingin mendapatkan **rank-k vector representation** dari sebuah kata di **vector space baru (rank-k subspace)**:

=> Vektor baris pada $U_k \times \Sigma_k$

Dari Contoh sebelumnya:

2-Dimensional Document Embedding dari “Cancer” adalah

Cancer	0.66	0.33
--------	------	------

 \times

Σ_2	18.93	0
	0	14.49

 $=$

Cancer	12.49	4.78
--------	-------	------


Jika ada Query, bagaimana hitung $\text{sim}(Q, D)$?

Vektor **Query** yang masih berada di dimensi awal perlu dipetakan ke **LSA (Semantic) Space** yang berukuran **K** dengan cara:

$$q_k = U_k^T \times q$$

Proof?

$$\text{sim}(q, d) = \text{sim}(q_k, d_k)$$



Vektor kolom pada $\Sigma_k \times V_k^T$
yang terasosiasi dengan **d**.

Uji Pemahaman Terhadap LSA

Refleksi

- Apa arti “distributional” pada “distributional word representations”?
- Apa itu “sparse word representations”? Apa contohnya?
- Apa itu “dense word representations”? Apa kelebihanannya?
- Apa itu term-document matrix? Apa saja yang dapat digunakan untuk mengisi setiap cell pada matrix tersebut?

Refleksi

Misal, kita melakukan **Singular Value Decomposition** terhadap term-document matrix C (berukuran $M \times N$):

$$C = U \times \Sigma \times V^T$$

- Informasi apa yang ada pada U ?
- Informasi apa yang ada pada Σ ? Apa makna singular values?
- Informasi apa yang ada pada V^T ?
- Bagaimana menghitung U , Σ , dan V^T ?
- Dimanakah *vector representation of words* berada?
- Bagaimana caranya mendapatkan *low-dimensional vector of words (dense)* berukuran $K (< \min(M, N))$?
- Bagaimana caranya mendapatkan *low-dimensional vector of documents (dense)* berukuran $K (< \min(M, N))$?

Refleksi

Perhatikan Term-Document matrix berikut (setiap cell berisi informasi TF):

	D1	D2	D3	D4	D5
Cancer	6	0	10	1	7
Flower	2	8	1	9	0
Tumor	6	2	7	0	8
Rose	1	6	0	7	1

Sebenarnya apa sih yang dilakukan LSA?

When forced to squeeze the terms/documents down to a k-dimensional space, the SVD should bring together terms with similar co-occurrences.

Refleksi

Perhatikan Term-Document matrix berikut (setiap cell berisi informasi TF):

Kira-kira ada berapa latent topics yang penting di sini?

	D1	D3	D5	D2	D4
Cancer	6	10	7	0	1
Tumor	6	7	8	2	0
Flower	2	1	0	8	9
Rose	1	0	1	6	7

LSA akan memindahkan term (baris) dan dokumen (kolom) sehingga vektor-vektor yang mirip akan berdekatan (ter-cluster).

Setelah ter-cluster, LSA kemudian bisa menemukan, kira-kira ada berapa “latent topics penting” yang ada pada koleksi dokumen tersebut.

	D1	D2	D3	D4	D5
Cancer	6	0	10	1	7
Flower	2	8	1	9	0
Tumor	6	2	7	0	8
Rose	1	6	0	7	1

Itulah mengapa hanya ada 2 singular values yang sangat besar dibandingkan 2 yang terkecil.

18.93	0	0	0
0	14.49	0	0
0	0	2.60	0
0	0	0	0.86

	1	2	3	4
Cancer	0.66	0.33	0.64	0.18
Flower	0.33	-0.71	0.18	-0.57
Tumor	0.61	0.25	-0.72	-0.19
Rose	0.24	-0.55	-0.18	0.77

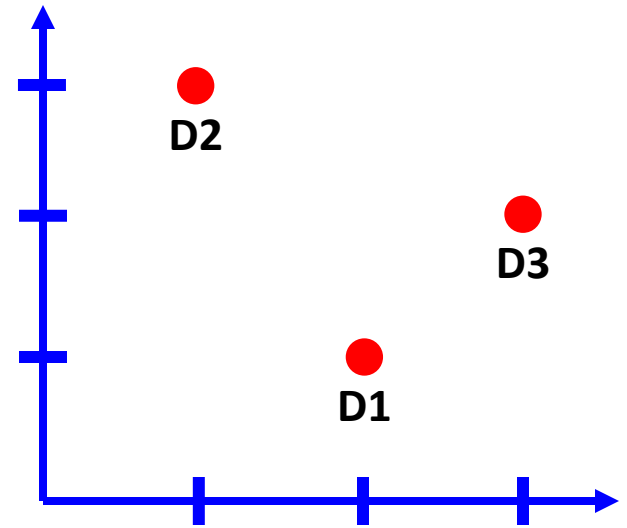
	D1	D2	D3	D4	D5
1	0.45	0.28	0.59	0.28	0.51
2	0.10	-0.59	0.30	-0.69	0.26
3	-0.11	-0.41	0.59	0.38	-0.56
4	-0.51	-0.42	-0.12	0.44	0.58
5	-0.70	0.46	0.42	-0.30	0.04

```
import numpy as np; u, s, vt = np.linalg.svd(C, full_matrices = True)
```

Simple Explanation of LSA

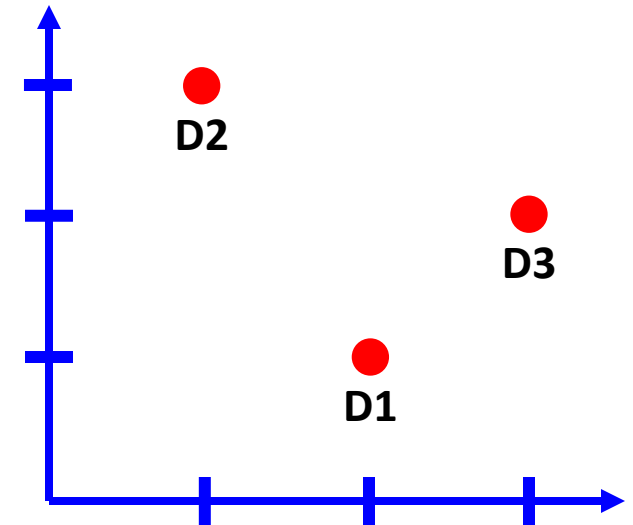
Misal, kita mempunyai term-document matrix berikut (bobot TF):

$$C = \begin{bmatrix} 2 & 1 & 3 \\ 1 & 3 & 2 \end{bmatrix}$$



Misal, kita mempunyai term-document matrix berikut (bobot TF):

$$C = \begin{bmatrix} 2 & 1 & 3 \\ 1 & 3 & 2 \end{bmatrix}$$



Misal, kita lakukan SVD:

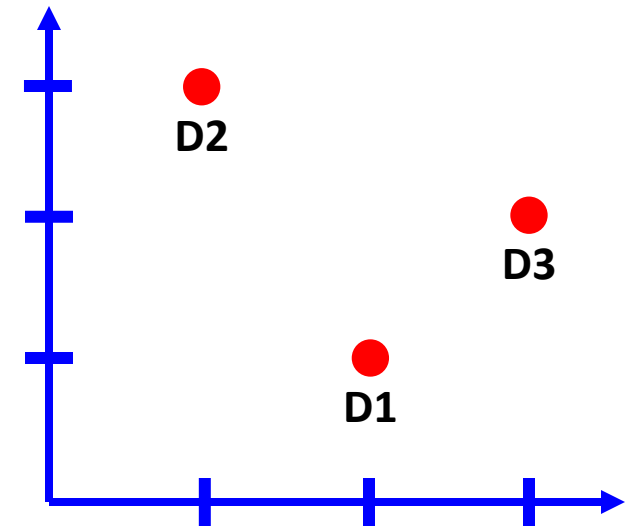
$$U = \begin{bmatrix} 0.71 & -0.71 \\ 0.71 & 0.71 \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} 5 & 0 & 0 \\ 0 & 1.7 & 0 \end{bmatrix}$$

$$V^T = \begin{bmatrix} 0.42 & 0.57 & 0.71 \\ -0.41 & 0.82 & -0.41 \\ -0.81 & -0.11 & 0.8 \end{bmatrix}$$

Misal, kita mempunyai term-document matrix berikut (bobot TF):

$$C = \begin{bmatrix} 2 & 1 & 3 \\ 1 & 3 & 2 \end{bmatrix}$$



Rank-1 approximation dari C:

$$U = \begin{bmatrix} 0.71 & -0.71 \\ 0.71 & 0.71 \end{bmatrix}$$

$$\Sigma_1 = \begin{bmatrix} 5 & 0 & 0 \\ 0 & \mathbf{0} & 0 \end{bmatrix}$$

$$V^T = \begin{bmatrix} 0.42 & 0.57 & 0.71 \\ -0.41 & 0.82 & -0.41 \\ -0.81 & -0.11 & 0.8 \end{bmatrix}$$

$$C_1 = U \times \Sigma_1 \times V^T = \begin{bmatrix} 1.5 & 2 & 2.5 \\ 1.5 & 2 & 2.5 \end{bmatrix}$$

Misal, kita mempunyai term-document matrix berikut (bobot TF):

$$C = \begin{bmatrix} 2 & 1 & 3 \\ 1 & 3 & 2 \end{bmatrix}$$



Rank-1 approximation dari C (Rank-1 Document Vector):

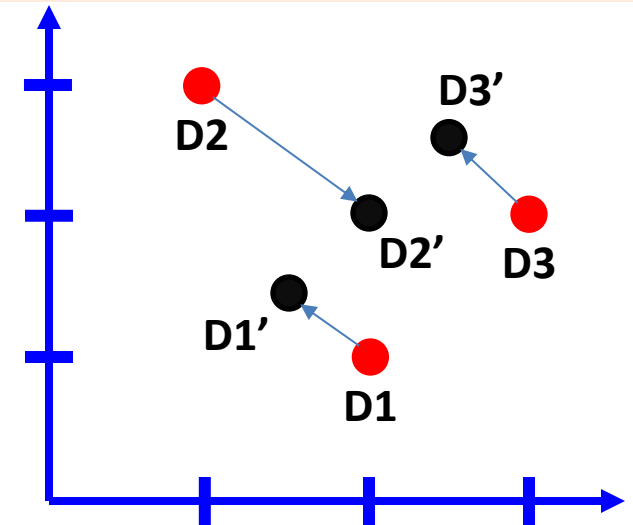
$$U = \begin{bmatrix} 0.71 & -0.71 \\ 0.71 & 0.71 \end{bmatrix}$$

$$\Sigma_1 = \begin{bmatrix} 5 & 0 & 0 \\ 0 & \mathbf{0} & 0 \end{bmatrix}$$

$$V^T = \begin{bmatrix} 0.42 & 0.57 & 0.71 \\ -0.41 & 0.82 & -0.41 \\ -0.81 & -0.11 & 0.8 \end{bmatrix}$$

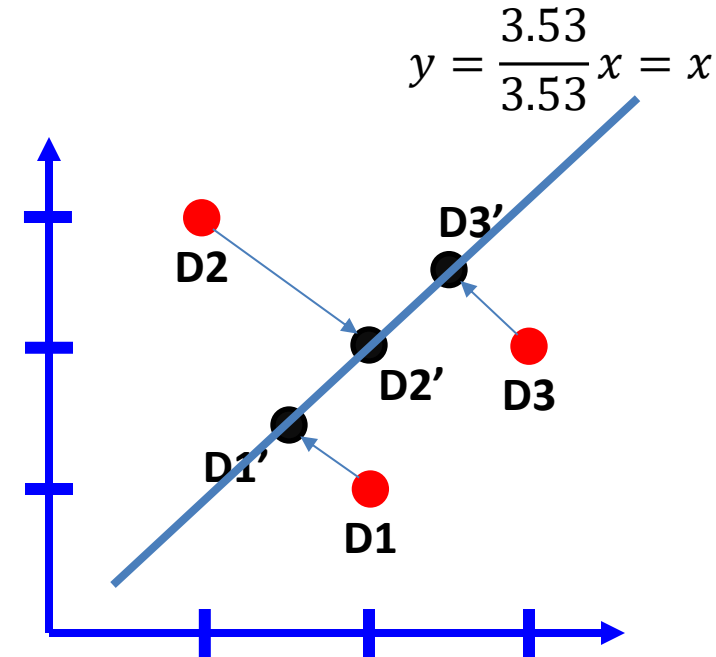
$$C_1 = U \times \Sigma_1 \times V^T = \begin{bmatrix} 1.5 & 2 & 2.5 \\ 1.5 & 2 & 2.5 \end{bmatrix}$$

We move the points to the **smallest squared distance**.



Misal, kita mempunyai term-document matrix berikut (bobot TF):

$$C = \begin{bmatrix} 2 & 1 & 3 \\ 1 & 3 & 2 \end{bmatrix}$$



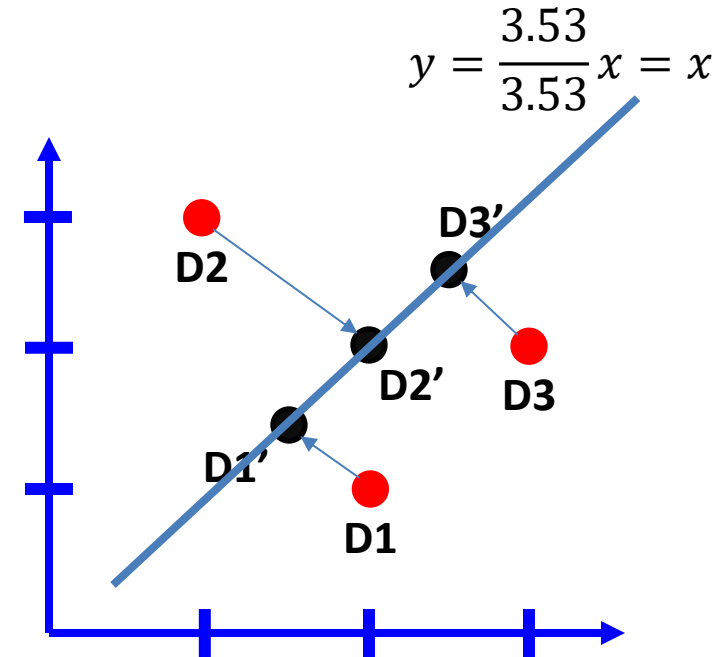
Vektor Kata:

$$U \times \Sigma_1 = \begin{bmatrix} 3.53 & 0 & 0 \\ 3.53 & 0 & 0 \end{bmatrix}$$

$$C_1 = U \times \Sigma_1 \times V^T = \begin{bmatrix} 1.5 & 2 & 2.5 \\ 1.5 & 2 & 2.5 \end{bmatrix}$$

Misal, kita mempunyai term-document matrix berikut (bobot TF):

$$C = \begin{bmatrix} 2 & 1 & 3 \\ 1 & 3 & 2 \end{bmatrix}$$



Vektor Dokumen di Space Dim = 1:

$$\Sigma_1 \times V^T = \begin{bmatrix} 2.1 & 2.8 & 3.5 \\ 0 & 0 & 0 \end{bmatrix}$$

$$C_1 = U \times \Sigma_1 \times V^T = \begin{bmatrix} 1.5 & 2 & 2.5 \\ 1.5 & 2 & 2.5 \end{bmatrix}$$

Misal, kita mempunyai term-document matrix berikut (bobot TF):

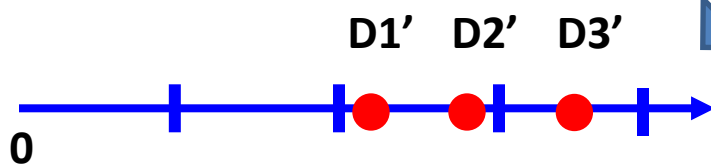
$$C = \begin{bmatrix} 2 & 1 & 3 \\ 1 & 3 & 2 \end{bmatrix}$$



Vektor Dokumen di Space Dim = 1:

$$\Sigma_1 \times V^T = \begin{bmatrix} 2.1 & 2.8 & 3.5 \\ 0 & 0 & 0 \end{bmatrix}$$

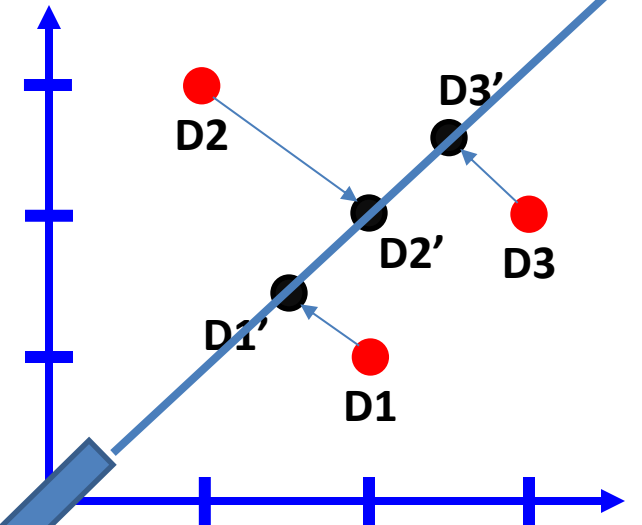
$$C_1 = U \times \Sigma_1 \times V^T = \begin{bmatrix} 1.5 & 2 & 2.5 \\ 1.5 & 2 & 2.5 \end{bmatrix}$$



Dimensi 1

Dimensi Original

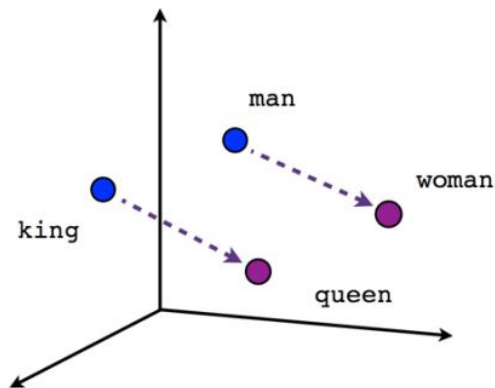
$$y = \frac{3.53}{3.53}x = x$$



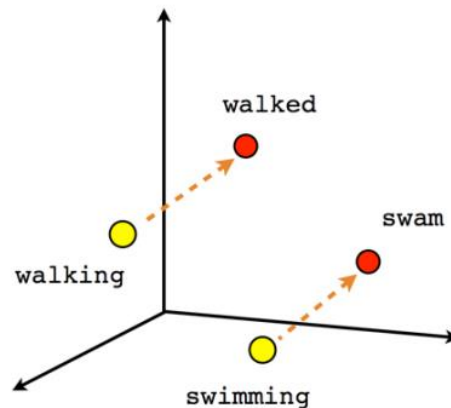
Neural Embeddings

Why Word Embeddings?

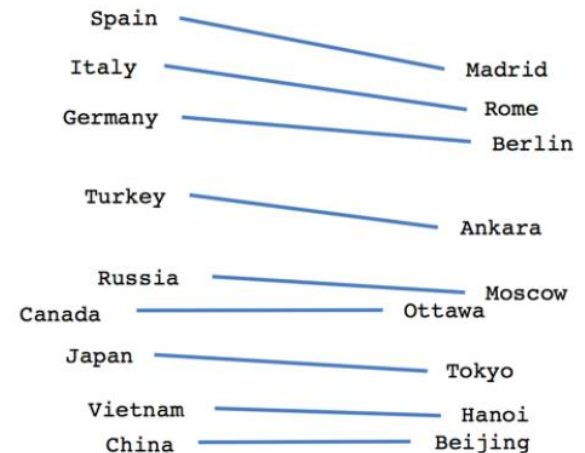
- Can capture the rich relational structure of the lexicon



Male-Female



Verb tense



Country-Capital

Word Embeddings

- Any technique that maps a word (or phrase) from **it's original high-dimensional sparse input** space to **a lower-dimensional dense vector space**.
- Vectors whose relative similarities correlate with semantic similarity
- Such vectors are used both as an end in itself (for computing similarities between terms), and as a representational basis for downstream NLP tasks, such as POS tagging, NER, text classification, etc.

Continuous Representation of Words

The Differences:

- In **information retrieval**, LSA and topic models use *documents as contexts*.
 - Capture semantic relatedness (“boat” and “water”)
- **Distributional semantic models** use *words as contexts* (more natural in linguistic perspective)
 - Capture semantic similarity (“boat” and “ship”)

Word Embedding



DSMs or Word Embeddings

- Count-based model
 - first collecting context vectors and then reweighting these vectors based on various criteria
- Predictive-based model (neural network)
 - vector weights are directly set to optimally predict the contexts in which the corresponding words tend to appear
 - Similar words occur in similar contexts, the system naturally learns to assign similar vectors to similar words.

Distributional Semantic Models

Other classification based on (Baroni et al., ACL 2014)

- Count-based models
 - Simple VSMs
 - Singular Value Decomposition (Golub & VanLoan, 1996)
 - Non-negative Matrix Factorization (Lee & Seung, 2000)
- Predictive-based models (**neural network**)
 - Self Organizing Map
 - Bengio et al's **Word Embedding (2003)**
 - Mikolov et al's **Word2Vec (2013)**

Word Analogy Task

- **Father** is to **Mother** as **King** is to _____ ?
- **Good** is to **Best** as **Smart** is to _____ ?
- **Indonesia** is to **Jakarta** as **Malaysia** is to _____ ?
- It turns out that the previous Word-Context based vector model is good for such analogy task.

$$\mathbf{V}_{\text{king}} - \mathbf{V}_{\text{father}} + \mathbf{V}_{\text{mother}} = \mathbf{V}_{\text{queen}}$$

Word2Vec (Mikolov et al., 2013)

Word2Vec

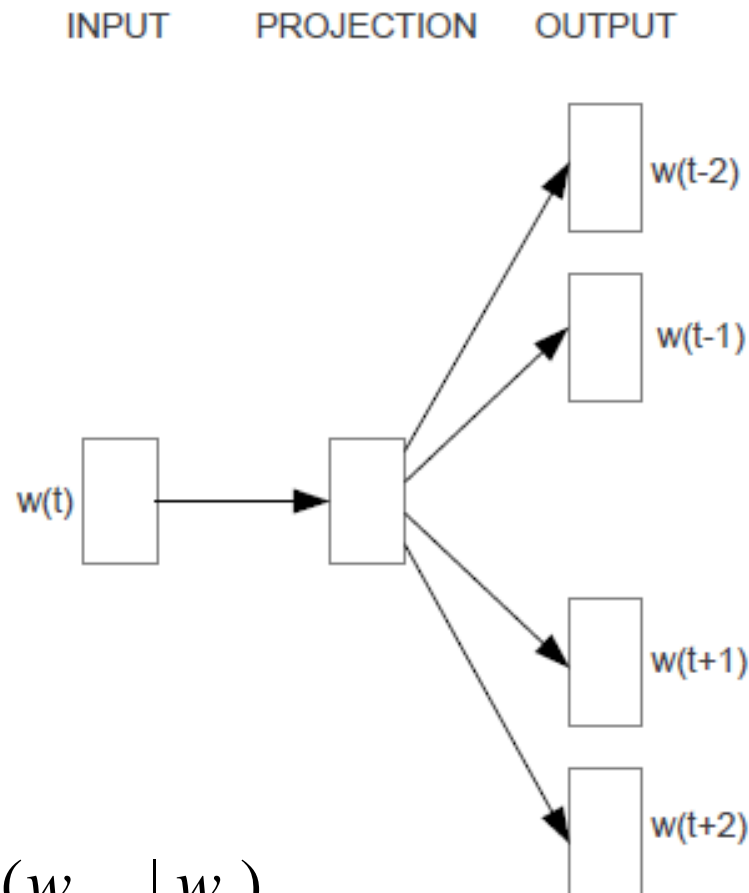
- One of the most popular Word Embedding models nowadays!
- There are two types of models:
 - Skip-Gram Model
 - Continuous Bag of Words Model (**CBOW**)

Skip-Gram

- **N-Gram language model** only looks at previous words as a context for predictions.
- This model tries to maximize classification of a word based on another word in the same sentence.
- Use each current word as an input to a log-linear classifier with continuous projection layer, and predict words within a certain range before and after the current word.

We seek a model for

$$P(w_{t+j} | w_t)$$



Skip-Gram

Feed-Forward Process

$$P(w_{t+j} | w_t) = \frac{\exp(y_{w_{t+j}})}{\sum_{i \in V} \exp(y_i)}$$

$$y_{w_t} = W \cdot x$$

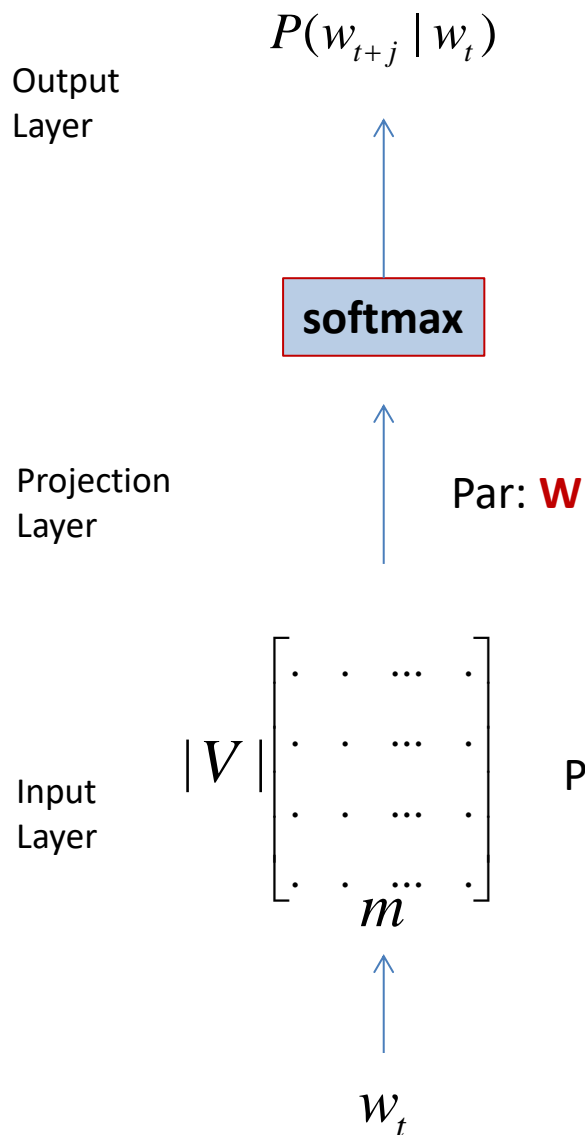
$$x = C(w_t)$$

Total Parameters:

$$\theta = \{W, C\}$$

$$C \in R^{|V| \times m}$$

$$W \in R^{m \times |V|}$$



Skip-Gram

Sudut Pandang Lain

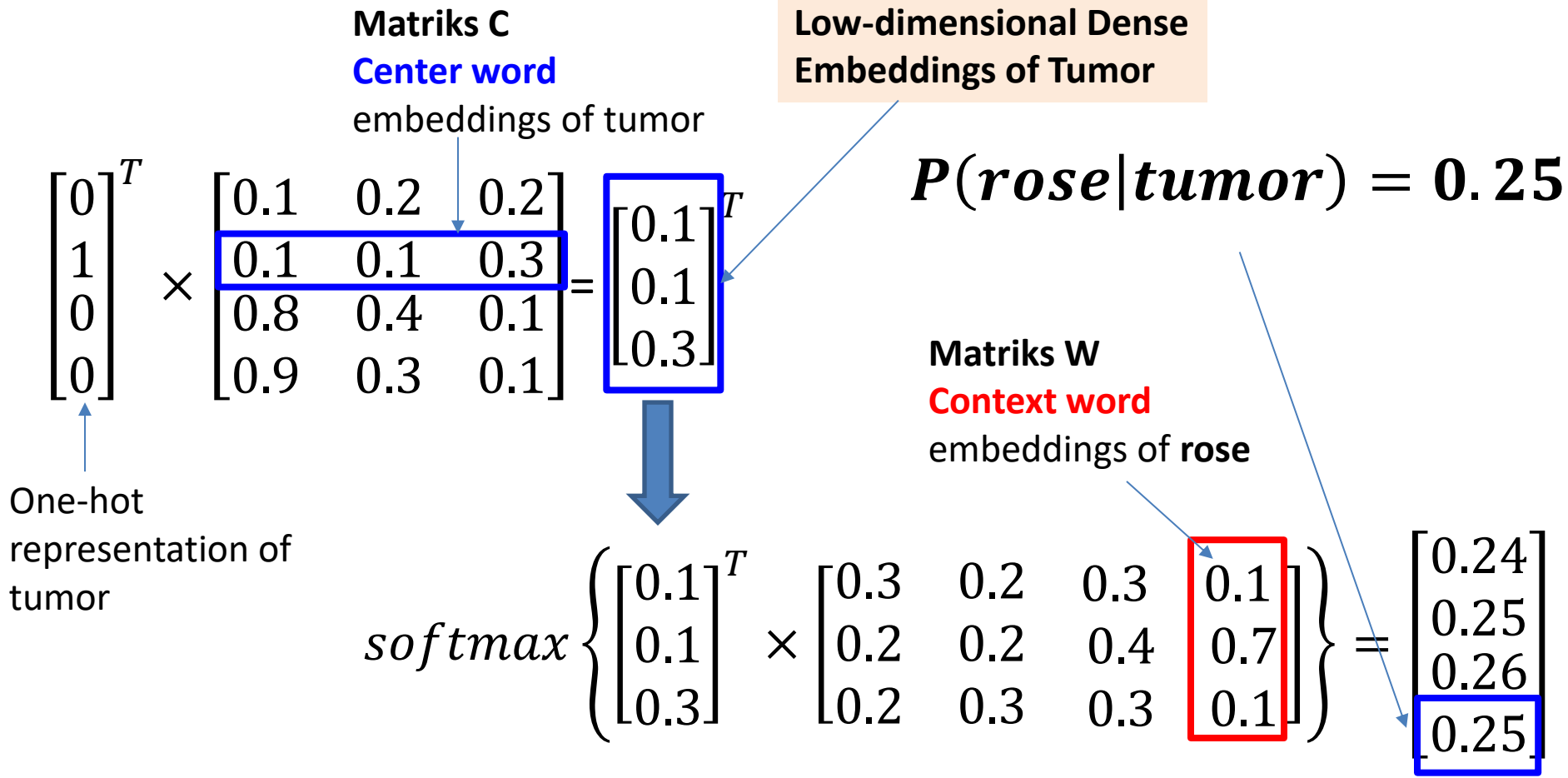
Sebuah kata \mathbf{t} terasosiasi dengan dua buah vector:

- Vector ketika \mathbf{t} berperan sebagai **center word** (baris pada matriks \mathbf{C})
- Vector ketika \mathbf{t} berperan sebagai **context word** (kolom pada matriks \mathbf{W})

Skip-Gram

Sudut Pandang Lain

Misal Vocab = [cancer, tumor, flower, rose]



Skip-Gram

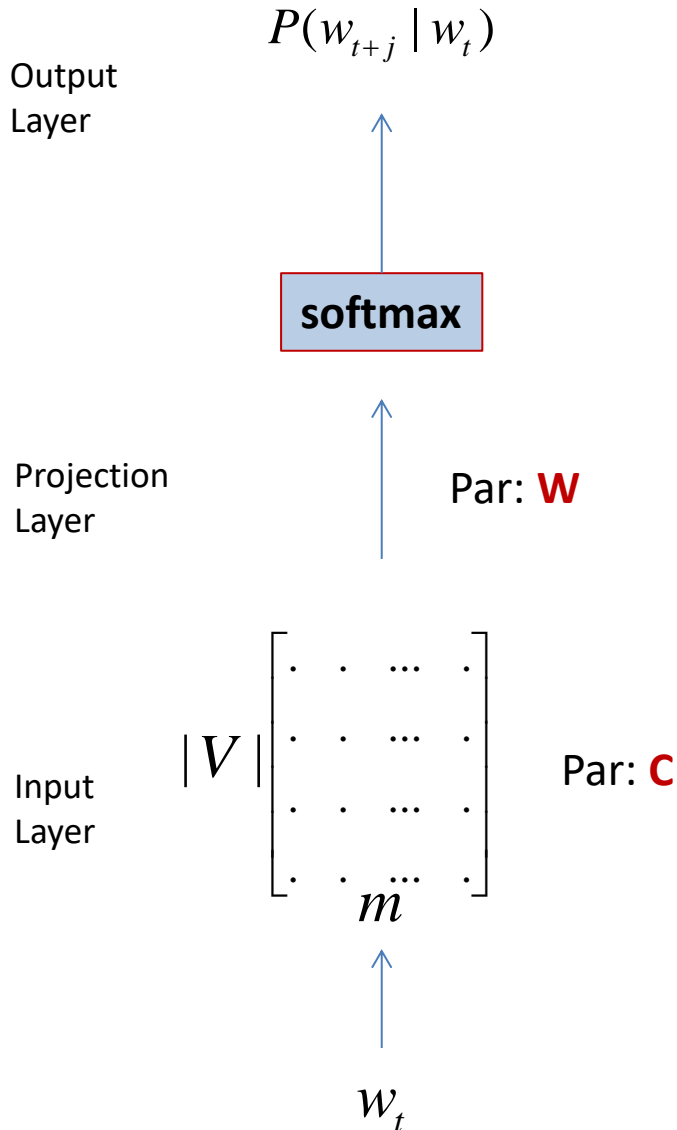
Dengan kata lain, Skip-Gram juga bisa dinyatakan dengan:

$$P(w_{t+j}|w_t) = \frac{\exp(s(w_{t+j}, w_t))}{\sum_{i \in V} \exp(s(w_i, w_t))}$$

$$s(w_i, w_t) = \text{center}(w_t)^T \cdot \text{context}(w_i)$$

Baris di matriks C

Kolom di matriks W

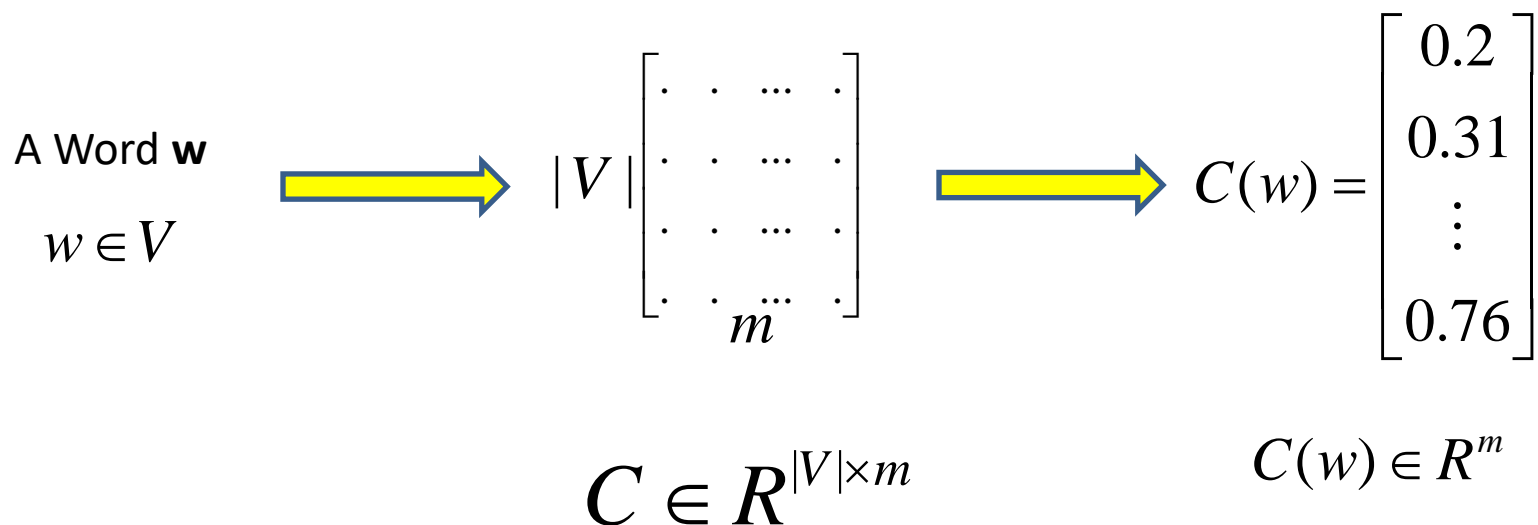


Where are the Word Embeddings?

- The previous model actually aims at building the language model.
 - So, **where is the Word Embedding model that we need?**
- **The answer is:** If you just need the Word Embedding model, you just need the matrix **C**

Where are the Word Embeddings?

After all parameters (including \mathbf{C}) are optimized, then we can use \mathbf{C} to map a word into its vector !



Skip-Gram

Training

Training is achieved by looking θ that maximizes the following Cost Function:

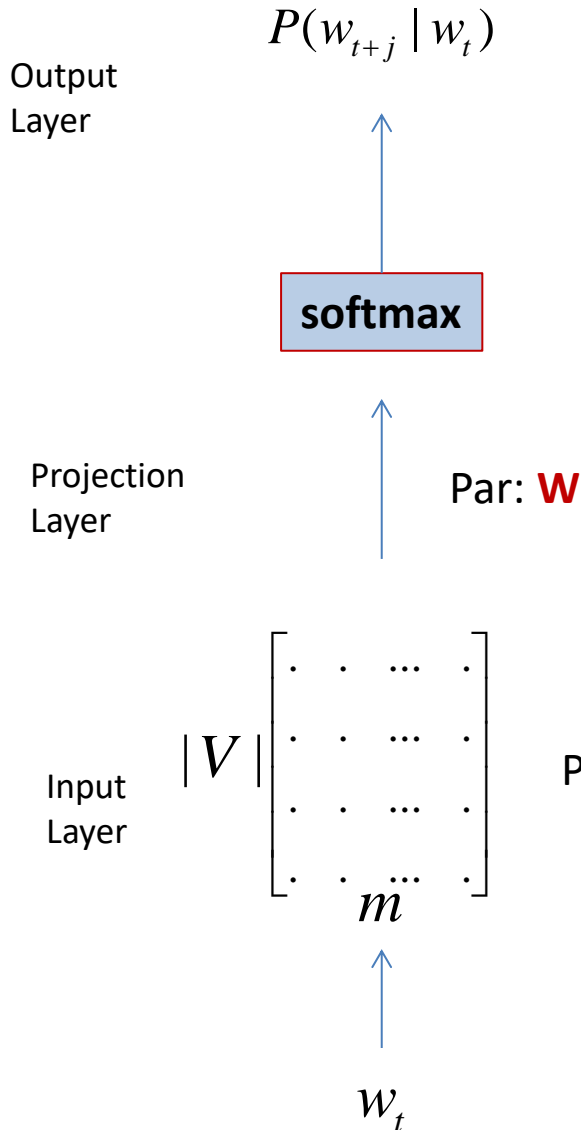
Given training data $w_1, w_2, w_3, \dots, w_{T-1}, w_T$

$$J(w_1 \dots w_T; \theta)$$

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

Regularization Terms

c is the maximum distance of the words, or **WINDOW size**



Skip-Gram

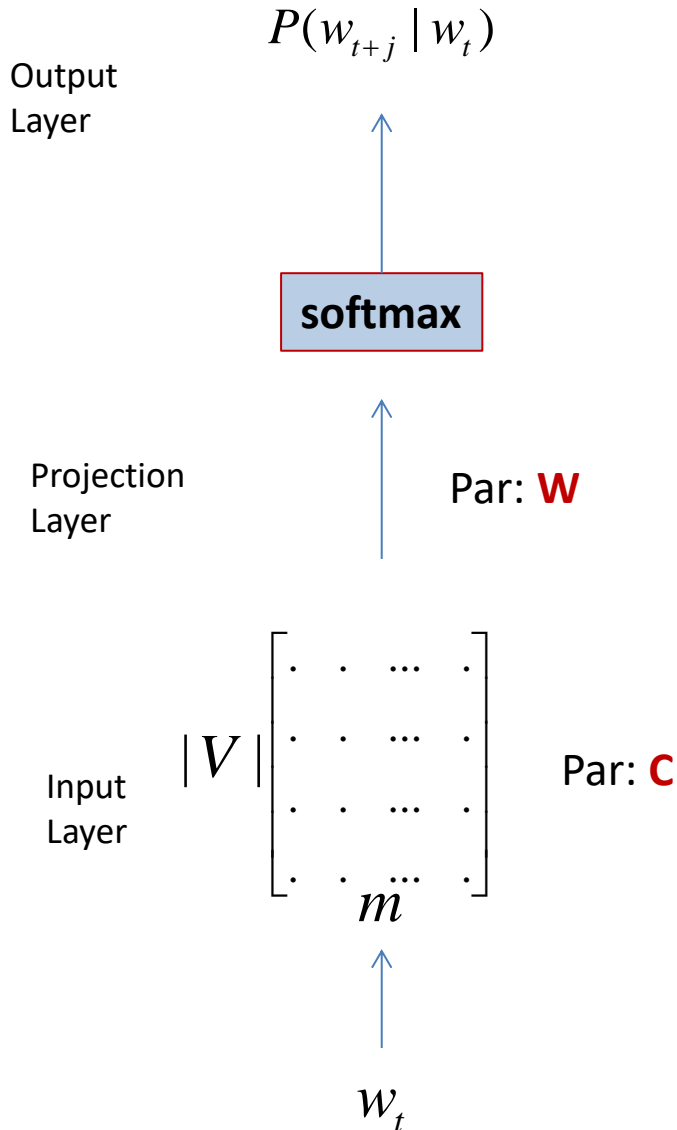
Training

Training is achieved by looking θ that maximizes the following Cost Function:

Given training data $w_1, w_2, w_3, \dots, w_{T-1}, w_T$

$$J(w_1 \dots w_T; \theta)$$

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



Proses optimasi (memaksimalkan) fungsi J = menggunakan **Categorical Cross Entropy** sebagai loss function pada ujung softmax layer.

Skip-Gram

How to develop dataset?

For example, let's consider the following dataset:

the quick brown fox jumped over the lazy dog

Using **c = 1 (or window = 1)**, we then have dataset:

([the, brown], quick), ([quick, fox], brown), ([brown, jumped], fox), ...

Therefore, our **(input, output)** dataset becomes:

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

For example, $P(w_{t+1} = \text{brown} \mid w_t = \text{quick})$

Use this dataset to learn $P(w_{t+j} \mid w_t)$

So that, the cost function is optimized!

$$J(w_1 \dots w_T; \theta) = \frac{1}{T} \sum_{t=1}^T \sum_{-c \leq j \leq c, j \neq 0} \log P(w_{t+j} \mid w_t)$$

Skip-Gram

Training

Actually, if we use **vanilla softmax**, then the computational complexity **per instance (Q)** is still costly.

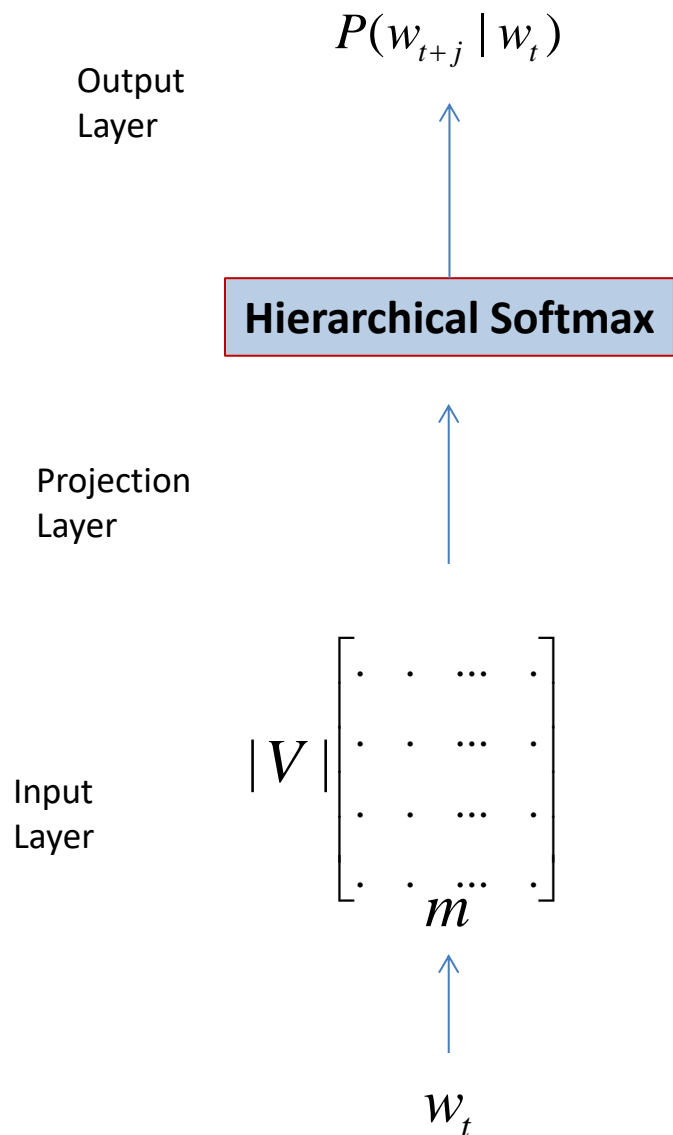
$$Q = D \times (m + m \times |V|)$$

D is the maximum distance of the words

To solve this problem, they use **Hierarchical Softmax layer**. This layer uses a **binary tree representation** of the output layer with $|V|$ units.

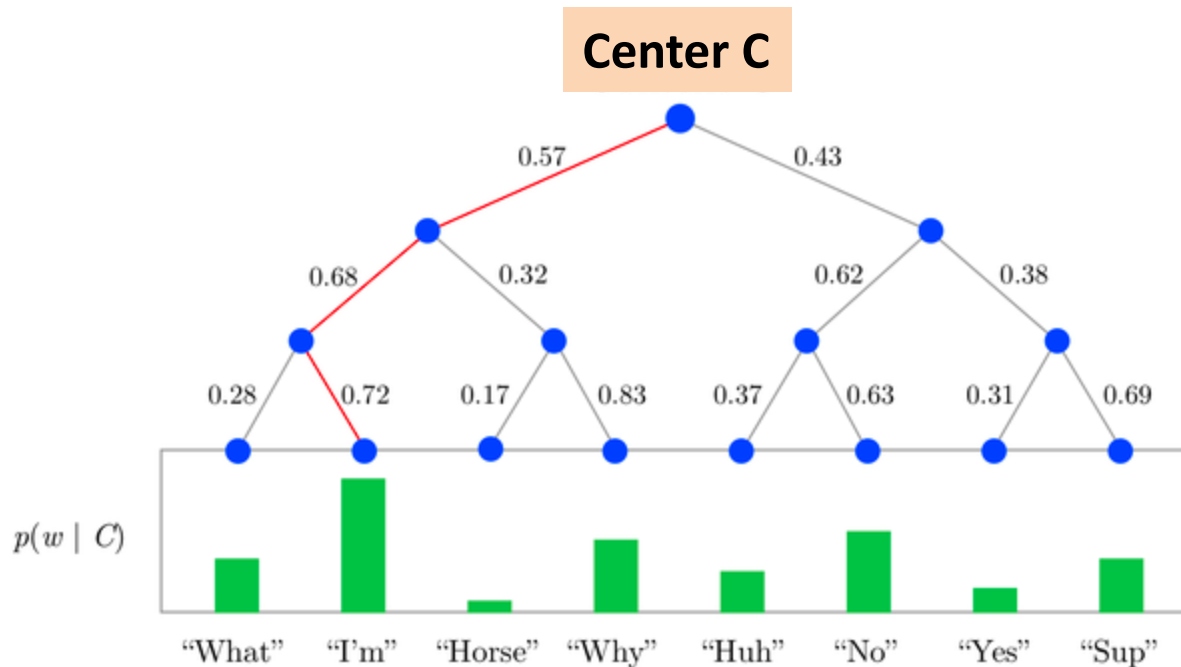
$$Q = D \times (m + m \times \log_2 |V|)$$

(Morin & Bengio, 2005)



Hierarchical Softmax

- A **multi-layer binary tree**
- The probability of a word is calculated through the **product of probabilities on each edge on the path to that node**.
- It is **$O(\log n)$** , compared to **$O(n)$** for vanilla softmax.



Hierarchical Softmax

- Misal, diberikan sebuah kata input “**kernel**” sebagai center, kita ingin memprediksi sebuah kata **konteks** **w**.
- Melakukan traversal Huffman Tree menggunakan kode binary yang diassign ke kata **w**.
- Perhitungan probability hanya melibatkan sekumpulan kecil node pada jalur root ke kata **w**.

Hierarchical Softmax

Huffman Tree

Meminimalkan expected search length

Gabung dua buah kata dengan frekuensi paling kecil; dan proses ini dilakukan terus.

kadal, 21

mutex, 15

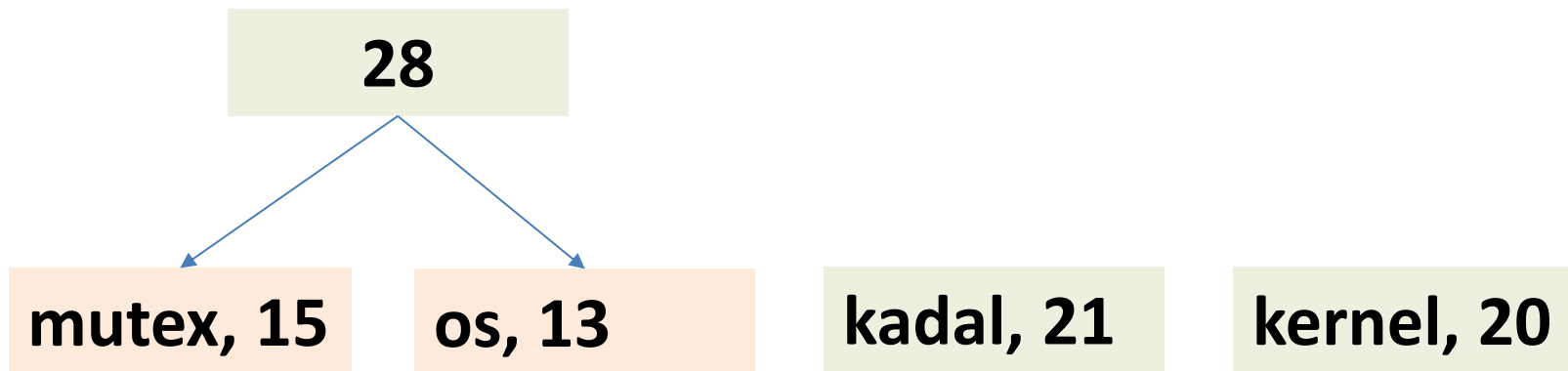
os, 13

kernel, 20

Hierarchical Softmax

Huffman Tree

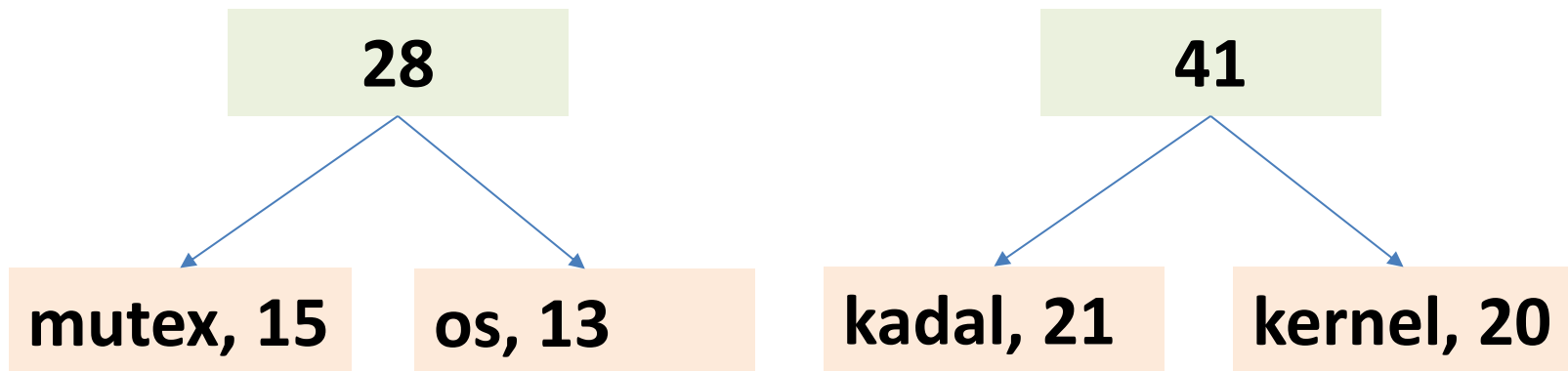
Gabung dua buah kata dengan frekuensi paling kecil; dan proses ini dilakukan terus.



Hierarchical Softmax

Huffman Tree

Gabung dua buah kata dengan frekuensi paling kecil; dan proses ini dilakukan terus.

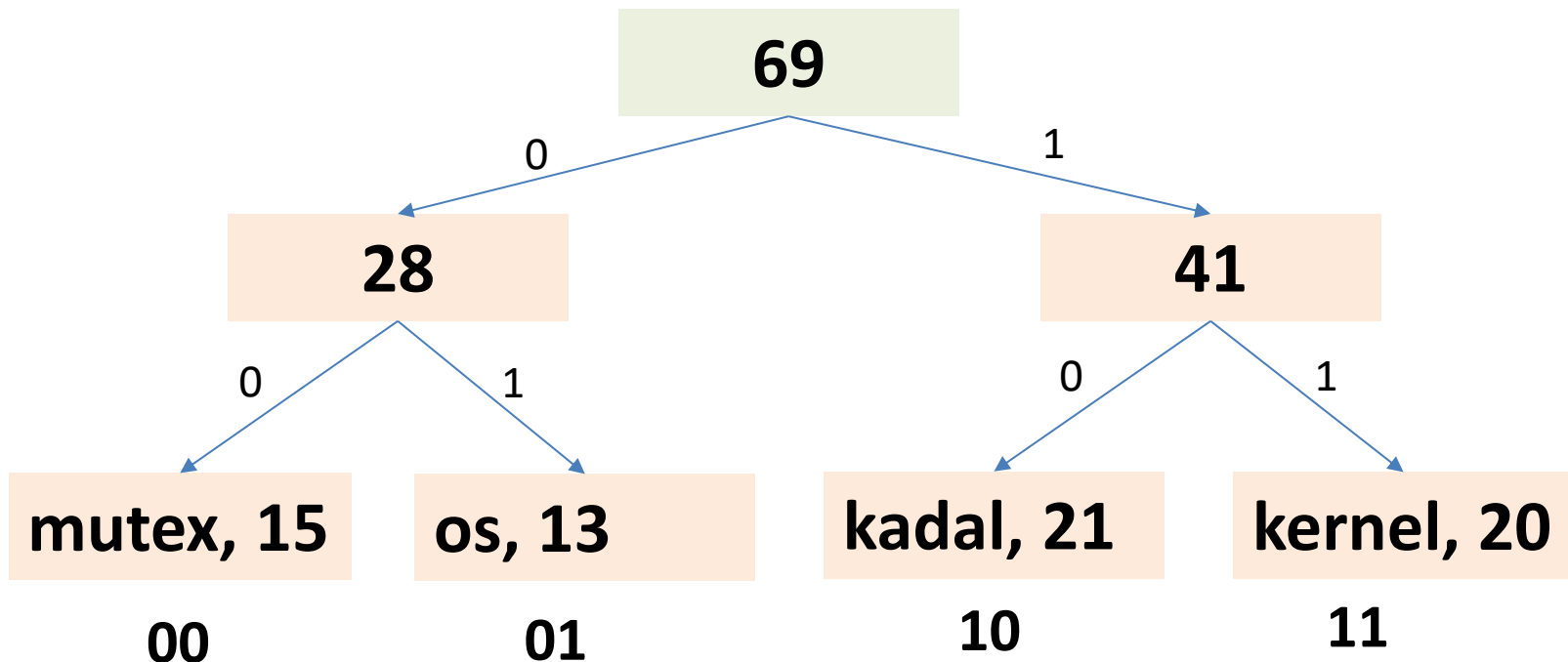


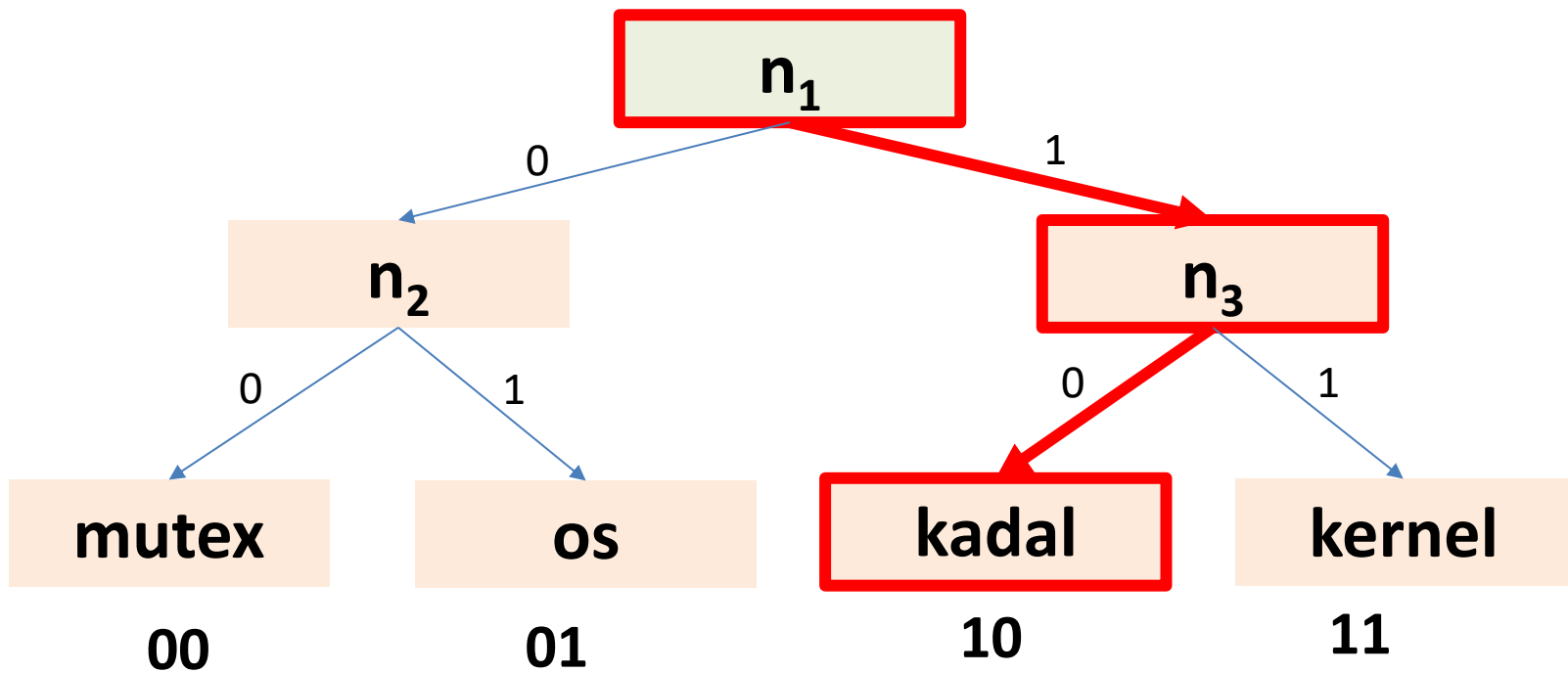
Hierarchical Softmax

Huffman Tree

Assign kode biner ke setiap kata di vocabulary

Gabung dua buah kata dengan frekuensi paling kecil; dan proses ini dilakukan terus.





$$\begin{aligned}
 P(kadal|c) &= P_{n_1}(\text{right}|c) \times P_{n_3}(\text{left}|c) \\
 &= (1 - P_{n_1}(\text{left}|c)) \times P_{n_3}(\text{left}|c)
 \end{aligned}$$

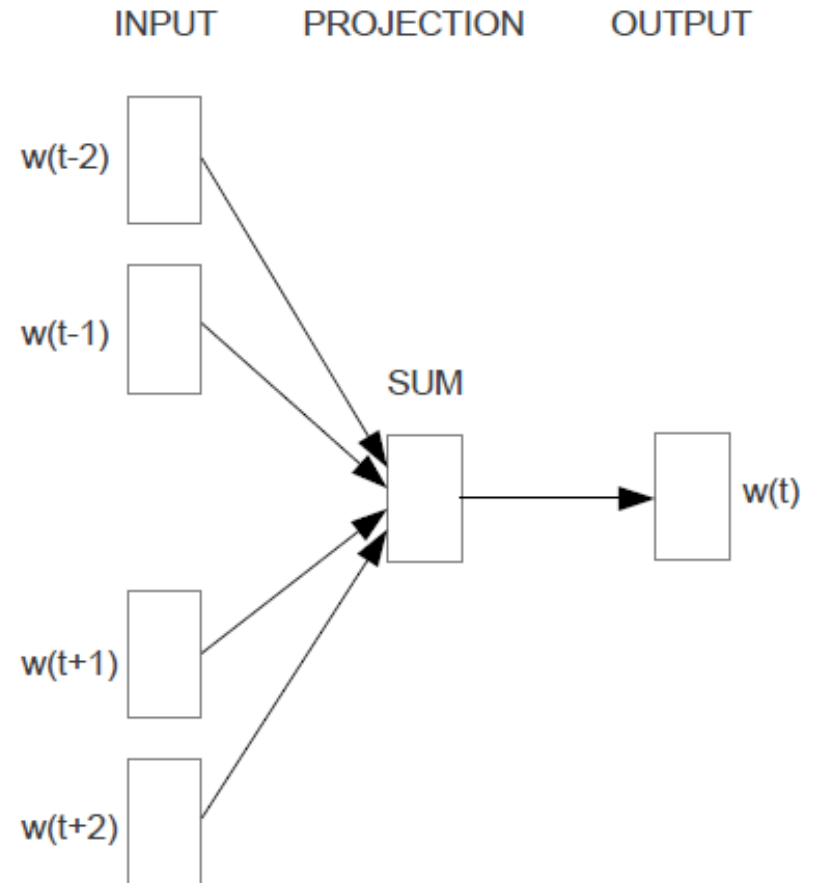
$$P_n(\text{left}|c) = \sigma(v_n^T \cdot \text{center}(c))$$

Trainable vector for node n

Hidden layer representation of the center c

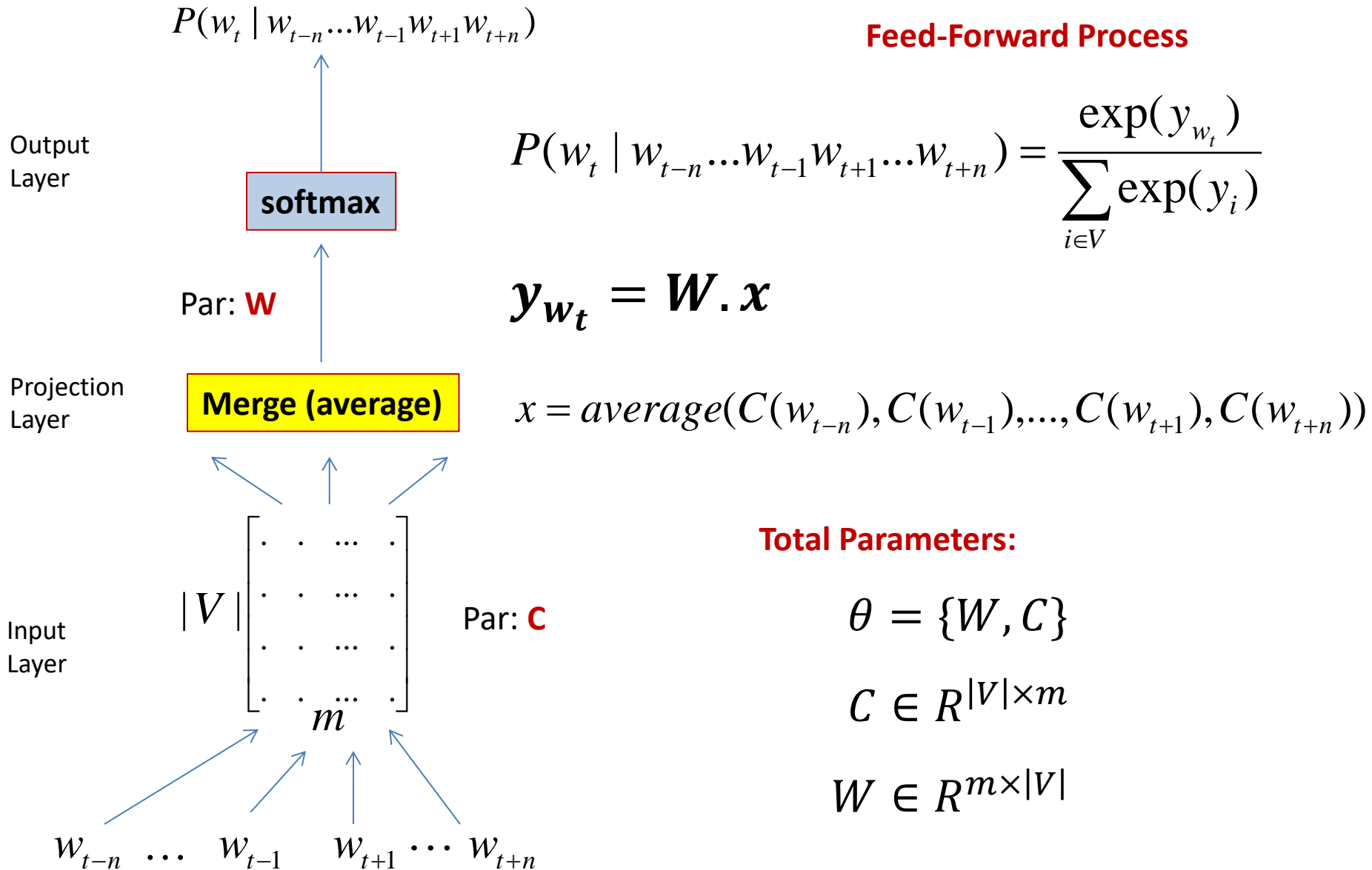
Continuous Bag-of-Words

- Mikolov's CBOW looks at **n words before and after the target words**.
 - Non-linear hidden layer is also removed.
 - All word vectors get projected into the same position (**their vectors are averaged**)
- “Bag-of-Words” is because **the order of words in the history does not influence the projection**.



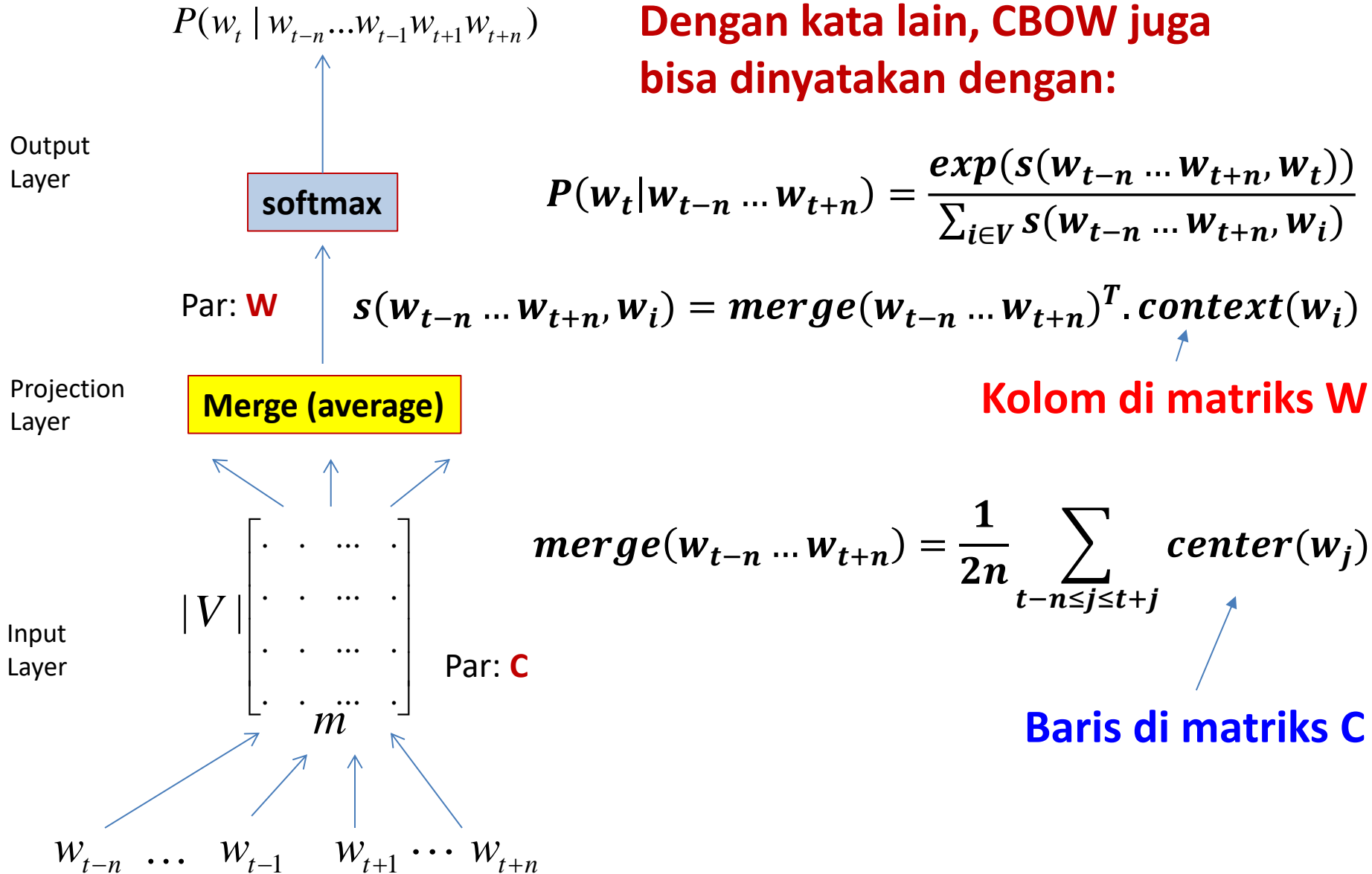
We seek a model for $P(w_t \mid w_{t-n} \dots w_{t-1} w_{t+1} \dots w_{t+n})$

Continuous Bag-of-Words

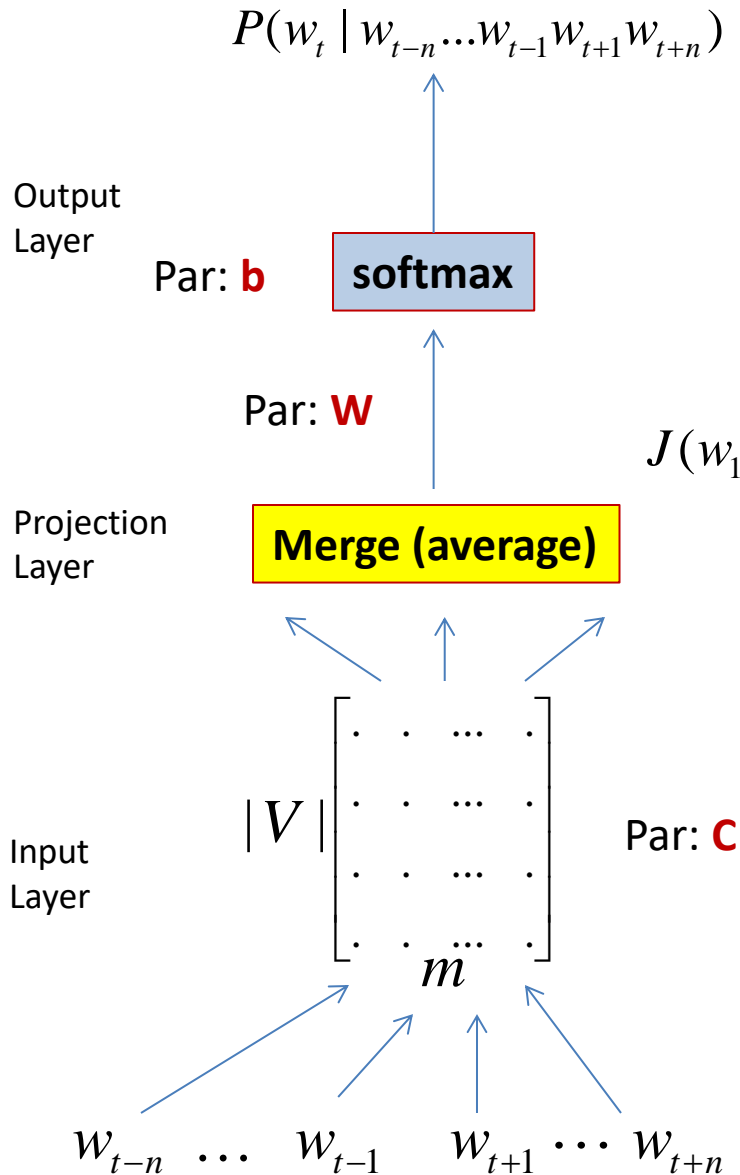


Continuous Bag-of-Words

Dengan kata lain, CBOW juga bisa dinyatakan dengan:



Continuous Bag-of-Words



Training

Training is achieved by looking θ that maximizes the following Cost Function:

Given training data $w_1, w_2, w_3, \dots, w_{T-1}, w_T$

$$J(w_1 \dots w_T; \theta) = \frac{1}{T} \sum_{t=1}^T \log P(w_t | w_{t-n} \dots w_{t-1} w_{t+1} \dots w_{t+n}) + R(\theta)$$

Regularization terms

Continuous Bag-of-Words

$$P(w_t | w_{t-n} \dots w_{t-1} w_{t+1} w_{t+n})$$

Training

Output
Layer

Hierarchical Softmax

Projection
Layer

Merge (average)

Input
Layer

$$|V| \begin{bmatrix} \cdot & \cdot & \dots & \cdot \\ \cdot & \cdot & \dots & \cdot \\ \cdot & \cdot & \dots & \cdot \\ \cdot & \cdot & \dots & \cdot \end{bmatrix}$$

m

$$w_{t-n} \dots w_{t-1} w_{t+1} \dots w_{t+n}$$

Actually, if we use **vanilla softmax**, then the computational complexity per instance (**Q**) is still costly.

$$Q = N \times m + m \times |V|$$

To solve this problem, they use **Hierarchical Softmax layer**. This layer uses a **binary tree representation** of the output layer with $|V|$ units.

$$Q = N \times m + m \times \log_2 |V|$$

(Morin & Bengio, 2005)

FastText (Bojanowski et al., TACL)

Subword Information

- Ekstensi dari Word2Vec -> bisa handle rare words atau kata yang out-of-vocabulary
- Perbedaannya adalah setiap kata tidak langsung diproyeksikan ke sebuah vector.
- Pada FastText, setiap kata “dipecah” dahulu menjadi “subwords”, dan setiap subword diproyeksikan ke vector.

Subword Information

- $w_t = \text{where}$
- Subwords dengan $\text{window} = 3$
 - <wh
 - whe
 - her
 - ere
 - re>
 - where



Kata asli itu sendiri diikutsertakan ke dalam set of ngrams

FastText – Skip-Gram Model

Given a word, what is the probability of its context word?

$$P(w_c | w_t) = \text{softmax}(s(w_c, w_t)) = \frac{\exp(s(w_c, w_t))}{\sum_{j=1}^{|Vocab|} \exp(s(w_j, w_t))}$$

Ingat bahwa, untuk kasus Word2Vec Skip-Gram biasa:

$$s(w_c, w_t) = \text{center}(w_t)^T \cdot \text{context}(w_c)$$

Untuk FastText, vektor kata adalah KOMBINASI dari vector Subwords !

Subword Information

- $\mathbf{w}_t = \text{where}$

Center embedding dari subword "<wh"

- Subwords dengan **window = 3**

– <wh $\rightarrow z_{g1} = [0.3, 0.1, \dots]$

– whe $\rightarrow z_{g2} = [0.1, 0.4, \dots]$

– her $\rightarrow z_{g3} = [0.5, 0.9, \dots]$

– ere $\rightarrow z_{g4} = [0.5, 0.5, \dots]$

– re> $\rightarrow z_{g5} = [0.3, 0.1, \dots]$

– **where** $\rightarrow z_{g6} = [0.1, 0.3, \dots]$

Kata asli itu sendiri diikutsertakan ke dalam set of ngrams

Subword Information

Representasi vector sebuah kata merupakan "sum" dari semua vector subwords.

- $\mathbf{w}_t = \text{where}$

- Subwords dengan $\text{window} = 3$

– <wh $\rightarrow z_{g1} = [0.3, 0.1, \dots]$

– whe $\rightarrow z_{g2} = [0.1, 0.4, \dots]$

– her $\rightarrow z_{g3} = [0.5, 0.9, \dots]$

– ere $\rightarrow z_{g4} = [0.5, 0.5, \dots]$

– re> $\rightarrow z_{g5} = [0.3, 0.1, \dots]$

– **where** $\rightarrow z_{g6} = [0.1, 0.3, \dots]$

$$s(w_c, w_t) = \sum_g z_g^T \cdot \text{context}(w_c)$$

Kata asli itu sendiri diikutsertakan ke dalam set of ngrams