

language model = cara yang assign probabilitas kepada sequence of words

word embedding = efek samping dari language model karena kita encode ke suatu baris

(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

Contoh: Path-based similarity

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

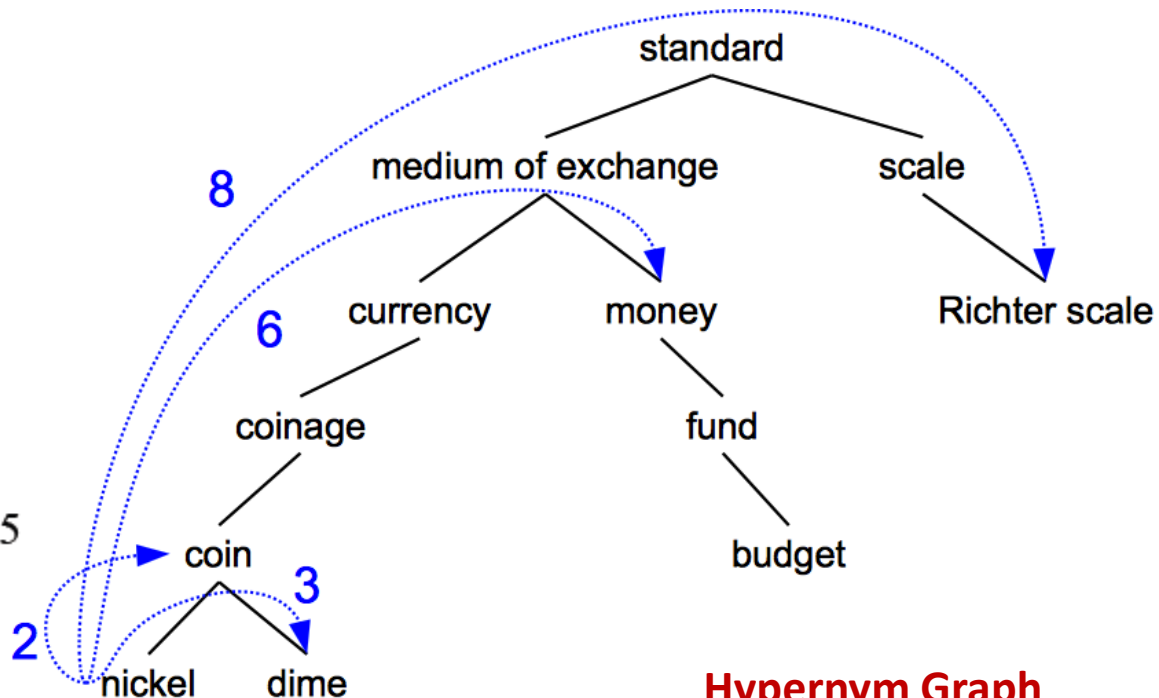
$$\text{simpath}(c_1, c_2) = \frac{1}{\text{pathlen}(c_1, c_2)}$$

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

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

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

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



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”

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.

	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.
- 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, $\text{rank}(C) \leq \min\{M, N\}$.

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

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

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

$$C\mathbf{x} = \lambda\mathbf{x}$$

are called the **eigenvalues** of C .

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

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

Contoh:

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

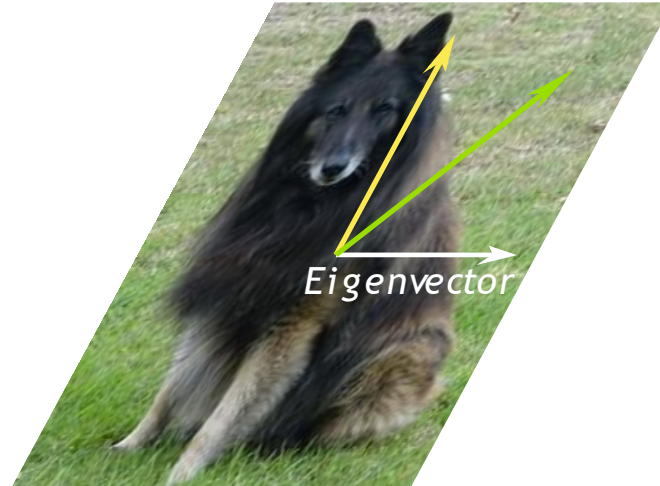
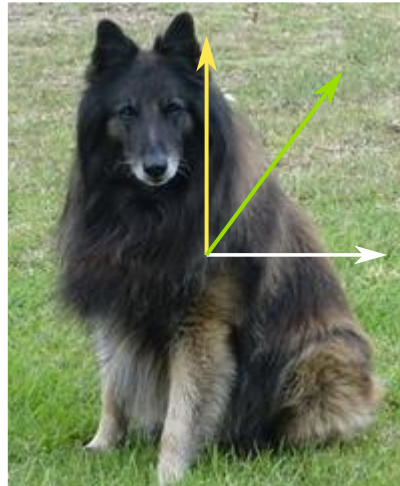
$$\mathbf{x}_1 = \begin{bmatrix} 1 \\ 4 \end{bmatrix} \quad \lambda_1 = 6$$

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

Alin Revisited: Eigenvector

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

“ditransformasi” = “dikali matriks”



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

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

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

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 \mathbf{C} adalah term-document matrix,

$\mathbf{C} \cdot \mathbf{C}^T$ mengandung dot products (similarity) dari semua pasangan **term vectors**; dan

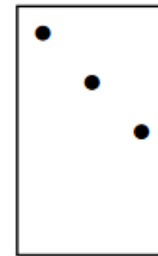
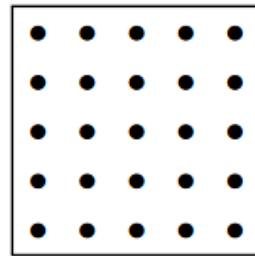
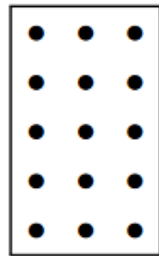
$\mathbf{C}^T \cdot \mathbf{C}$ mengandung dot products (similarity) dari semua pasangan **document vectors**.

Dekomposisi Term-Document Matrix!

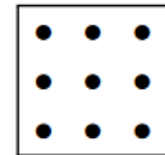
Misal, C adalah term-document matrix. C didekomposisi menjadi perkalian 3 matrix (**Singular Value Decomposition**).

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

Jika $M > N$



Cell selain titik = 0



C

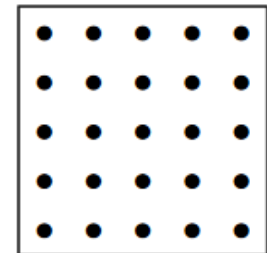
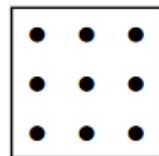
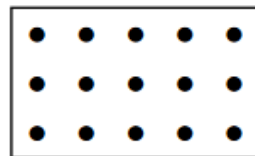
$=$

U

Σ

V^T

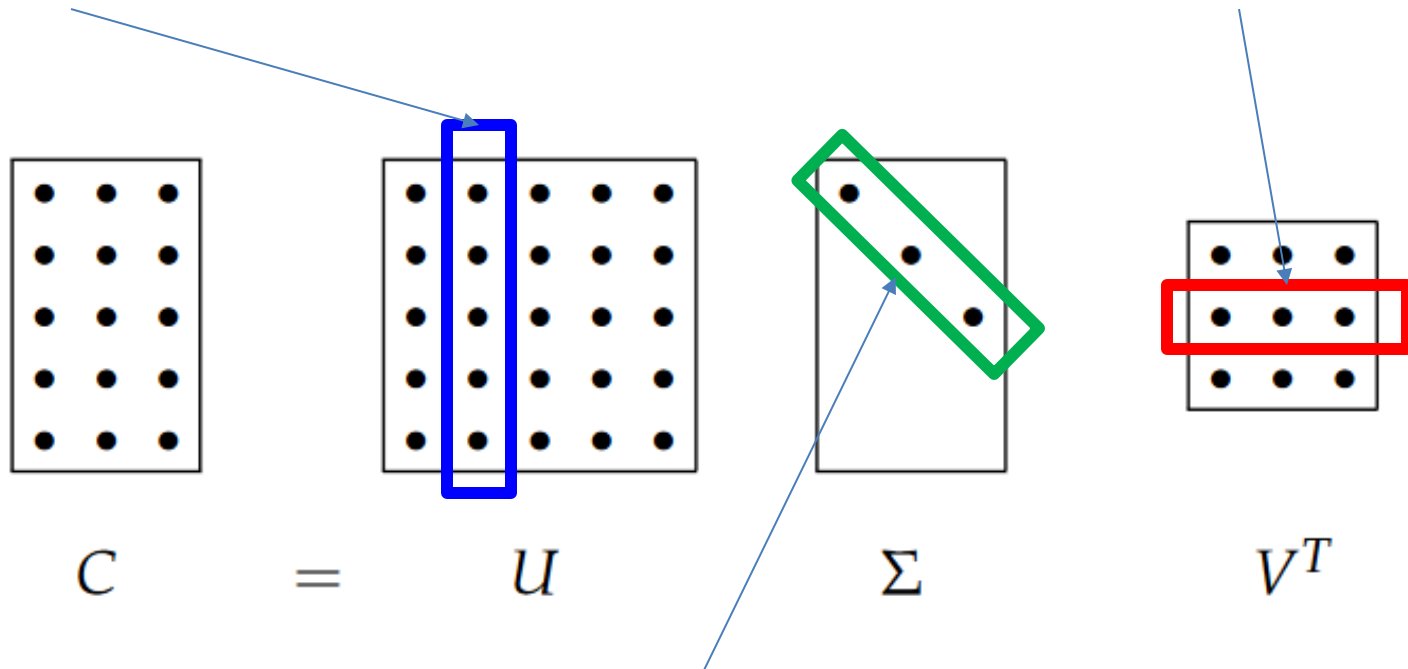
Jika $M < N$



Cara menghitung U , Σ , dan V^T

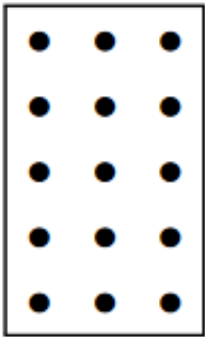
Kolom-kolom pada U
adalah **orthogonal
eigenvector** dari $C.C^T$

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



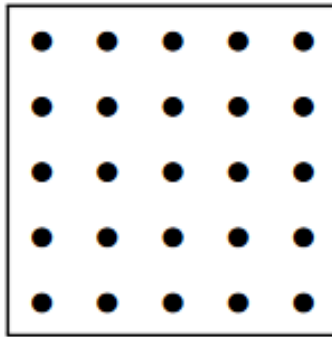
Singular Values (**Akar kuadrat dari Eigenvalues**)
dari $C^T.C$ atau $C.C^T$ (sama saja).

Truncated SVD

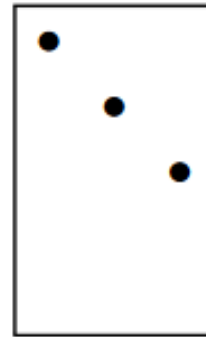


C

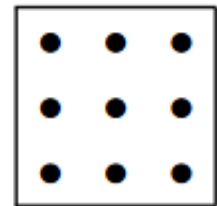
$=$



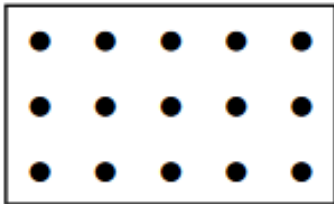
U



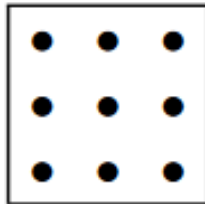
Σ



V^T



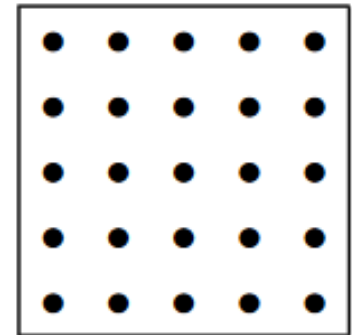
$M \times N$



$M \times M$



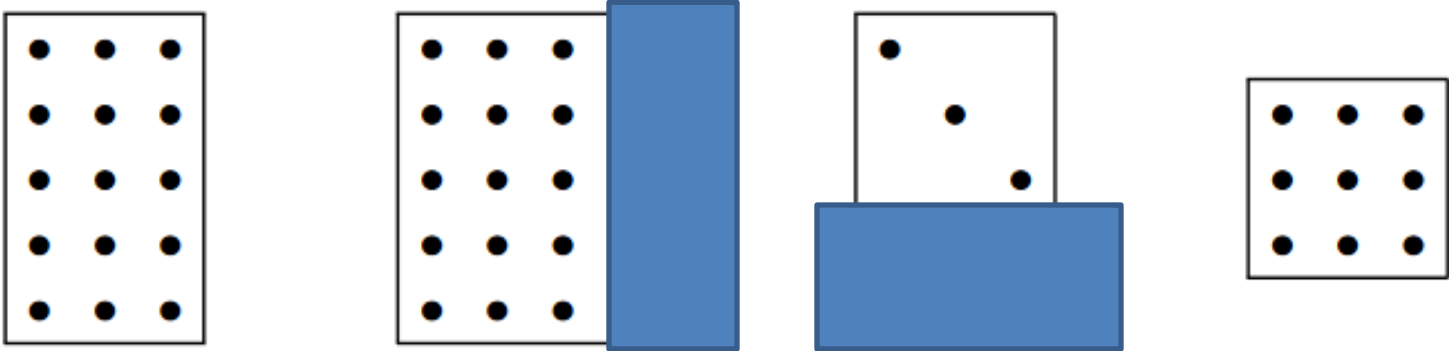
$M \times N$



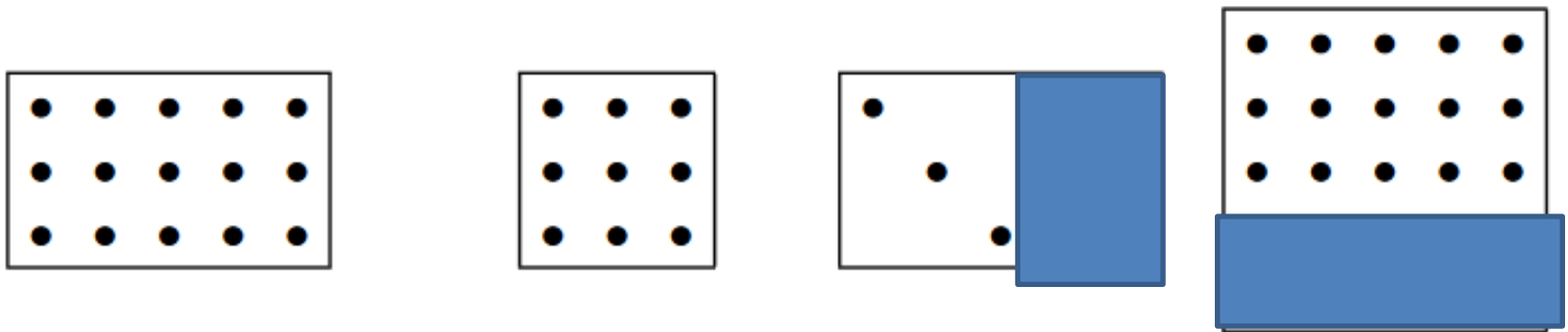
$N \times N$

Truncated SVD

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



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



$M \times N$

$M \times L$

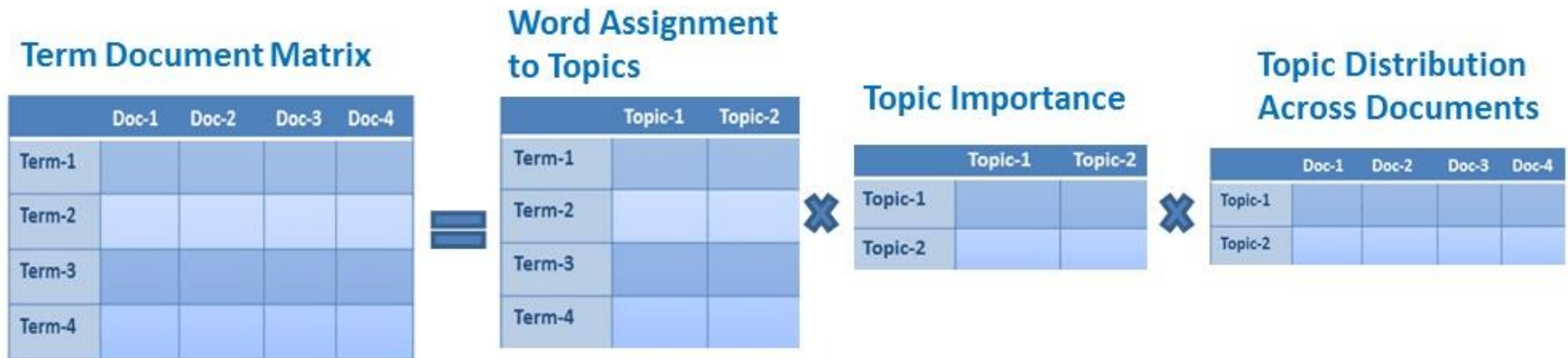
$L \times L$

$L \times N$

Interpretasi Truncated SVD

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

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



	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	truncated 0
0	14.49	0	0	0
0	0	2.60	0	0
0	0	0	0.86	0

	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

V^T

← truncated

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

C	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

Beberapa fakta:

Kolom-kolom U saling **orthogonal** dan merupakan **basis baru** untuk **dokumen-dokumen (kolom pada C)** di ruang yang baru (hasil rotasi).

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

truncated

	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

L2-Norm = 1

← truncated

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

C	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

Beberapa fakta:

Baris-baris V^T saling **orthogonal** dan merupakan **basis baru** untuk **term-term (baris pada C)** di ruang yang baru (hasil rotasi).

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

	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

V^T

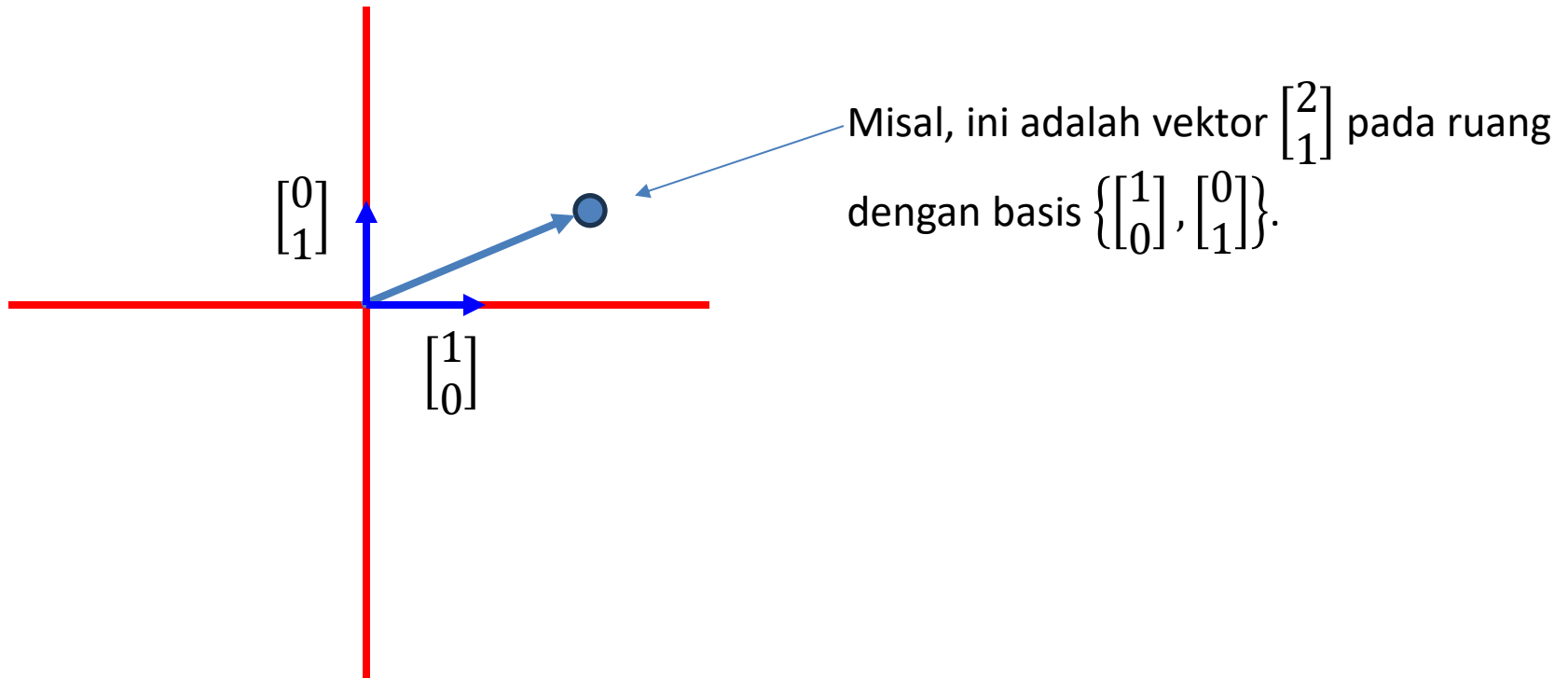
L2-Norm = 1

← truncated

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

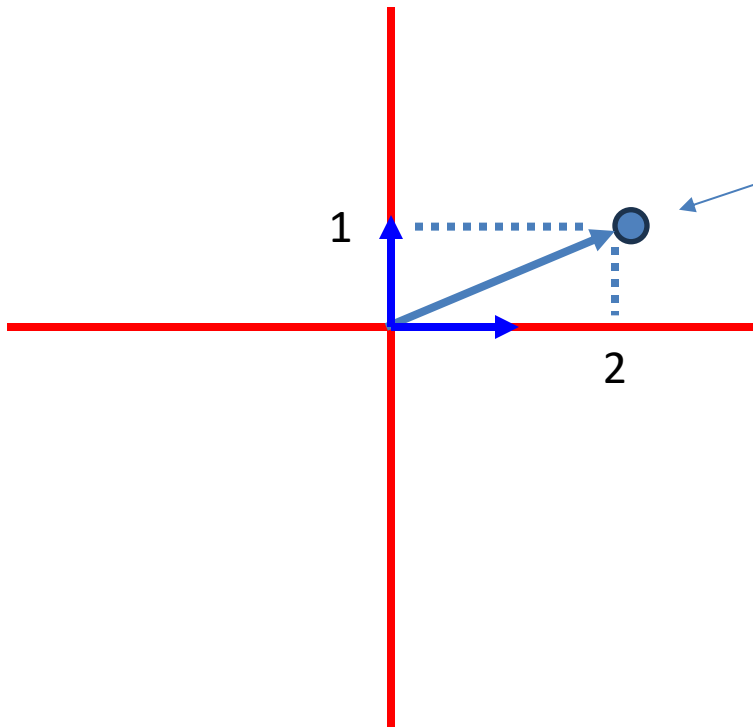

Review Aljabar Linear

Dot product dapat digunakan untuk memproyeksikan sebuah vektor ke himpunan basis-basis yang **orthogonal** pada sebuah ruang vektor (untuk dapat koordinat).



Review Aljabar Linear

Dot product dapat digunakan untuk memproyeksikan sebuah vektor ke himpunan basis-basis yang **orthogonal** pada sebuah ruang vektor (untuk dapat koordinat).



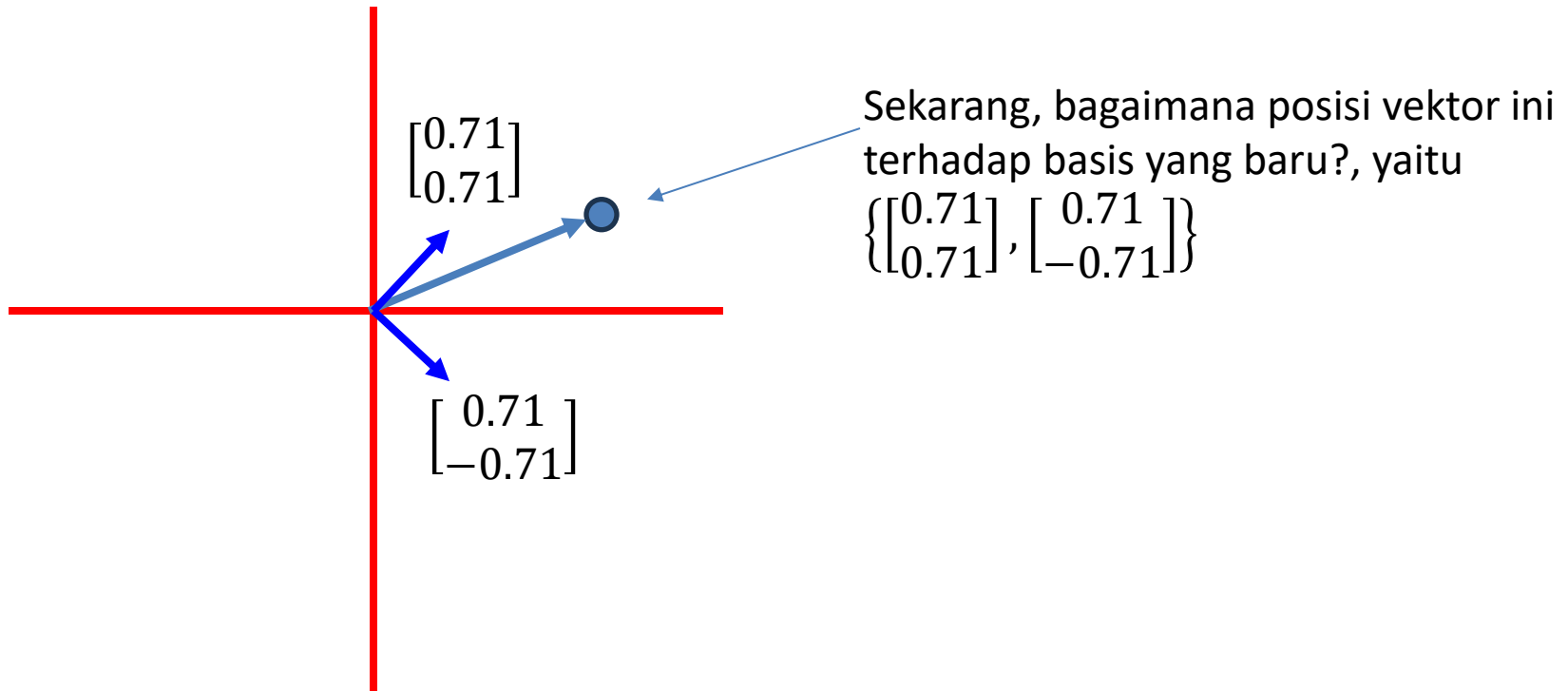
Misal, ini adalah vektor $\begin{bmatrix} 2 \\ 1 \end{bmatrix}$ pada ruang dengan basis $\left\{ \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \end{bmatrix} \right\}$.

$$\begin{bmatrix} 2 \\ 1 \end{bmatrix} \circ \begin{bmatrix} 1 \\ 0 \end{bmatrix} = 2 \quad \text{Koodinat di sumbu } \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} 2 \\ 1 \end{bmatrix} \circ \begin{bmatrix} 0 \\ 1 \end{bmatrix} = 1 \quad \text{Koodinat di sumbu } \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

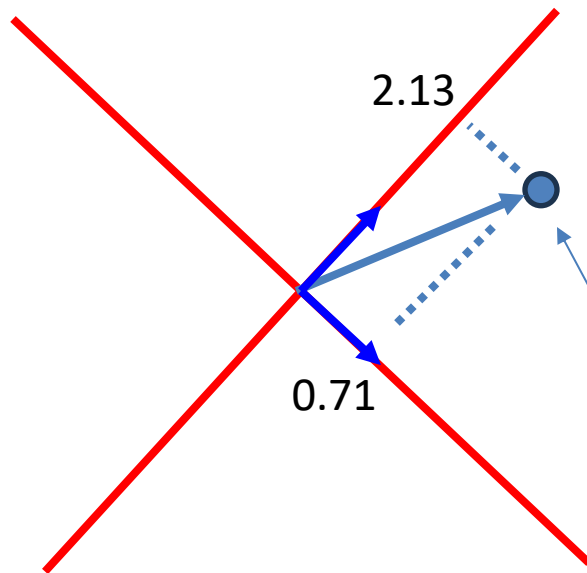
Review Aljabar Linear

Dot product dapat digunakan untuk memproyeksikan sebuah vektor ke himpunan basis-basis yang **orthogonal** pada sebuah ruang vektor (untuk dapat koordinat).



Review Aljabar Linear

Dot product dapat digunakan untuk memproyeksikan sebuah vektor ke himpunan basis-basis yang **orthogonal** pada sebuah ruang vektor (untuk dapat koordinat).



Sekarang, bagaimana posisi vektor ini terhadap basis yang baru, yaitu $\left\{ \begin{bmatrix} 0.71 \\ 0.71 \end{bmatrix}, \begin{bmatrix} 0.71 \\ -0.71 \end{bmatrix} \right\}$

$$\begin{bmatrix} 2 \\ 1 \end{bmatrix} \circ \begin{bmatrix} 0.71 \\ 0.71 \end{bmatrix} = 2.13$$

$$\begin{bmatrix} 2 \\ 1 \end{bmatrix} \circ \begin{bmatrix} 0.71 \\ -0.71 \end{bmatrix} = 0.71$$

ini adalah vektor $\begin{bmatrix} 2.13 \\ 0.71 \end{bmatrix}$ di ruang yang baru ini

C	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

Kolom-kolom **U** saling **orthogonal** dan merupakan **basis baru** untuk **dokumen-dokumen (kolom pada C)** di ruang yang baru (hasil rotasi).

Jika **ini** adalah Vektor Dokumen **D2** di ruang vektor original, apa Vektor Dokumen **D2** di ruang vektor yang baru dengan basis adalah kolom-kolom pada **U**?

	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

C	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		

Kolom-kolom **U** saling **orthogonal** dan merupakan **basis baru** untuk **dokumen-dokumen (kolom pada C)** di ruang yang baru (hasil rotasi).

DOT PRODUCT!

Jika **ini** adalah Vektor Dokumen **D2** di ruang vektor original, apa Vektor Dokumen **D2** di ruang vektor yang baru dengan basis adalah kolom-kolom pada **U**?

	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

Seperti yang sudah dipelajari pada review aljabar linear sebelumnya, ya sudah, Anda hanya perlu **dot-product** vektor original **D2** ke masing-masing 4 buah vektor basis pada kolom **U**.

C	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		

Kolom-kolom U saling **orthogonal** dan merupakan **basis baru** untuk **dokumen-dokumen (kolom pada C)** di ruang yang baru (hasil rotasi).

DOT PRODUCT!

Jika **ini** adalah Vektor Dokumen **D2** di ruang vektor original, apa Vektor Dokumen **D2** di ruang vektor yang baru dengan basis adalah kolom-kolom pada U ?

	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

Seperti yang sudah dipelajari pada review aljabar linear sebelumnya, ya sudah, Anda hanya perlu **dot-product** vektor original **D2** ke masing-masing 4 buah vektor basis pada kolom U ; **secara aljabar, vektor dokumen pada dimensi baru adalah kolom-kolom pada:**

$$U^T C$$

C	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		

Kolom-kolom **U** saling **orthogonal** dan merupakan **basis baru** untuk **dokumen-dokumen (kolom pada C)** di ruang yang baru (hasil rotasi).

DOT PRODUCT!

Jika **ini** adalah Vektor Dokumen **D2** di ruang vektor original, apa Vektor Dokumen **D2** di ruang vektor yang baru dengan basis adalah kolom-kolom pada **U**?

	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

Seperti yang sudah dipelajari pada review aljabar linear sebelumnya, ya sudah, Anda hanya perlu **dot-product** vektor original **D2** ke masing-masing 4 buah vektor basis pada kolom **U**; **secara aljabar, vektor dokumen pada dimensi baru adalah kolom-kolom pada:**

Kolomnya sekarang itu ada document embedding

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

Sama saja dengan kolom-kolom pada V^T , yang kemudian dikali bobot singular value pada topik yang bersesuaian.

	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

Apakah kolom pada V^T juga merupakan **Document Embedding**?

Kolom pada V^T sebenarnya bisa juga “dianggap” sebagai vektor “baru” dari sebuah dokumen di ruang vektor baru; tapi ini adalah versi yang **UNWEIGHTED**, atau versi yang belum dikali dengan bobot topik Σ .

	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

document embedding yang weighted ada di $U^T C$ atau $\Sigma * V^T$

Tapi kalau document embedding yang UNWEIGHTED ada di V^T

LSA: Low-rank approximation of C

Ada juga yang sebut LSI (Latent Semantic Indexing)

LSA (Latent Semantic Analysis)

Truncating U , Σ , V^T to K dimensions produces best possible K rank approximation of original matrix $C \rightarrow$ **rank K approximation to C with the smallest error (Frobenius Norm)**

Nah tujuannya kita ingin mencari aproksimasi C yang misalnya Ranknya terbaik (yang paling dekat sama C)

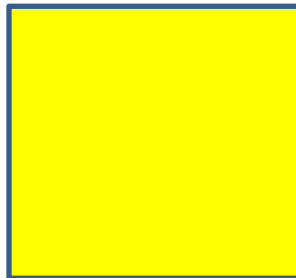
Buang (jadikan nol) $L - K$ singular values paling kecil di Σ !

$$U, M \times L$$



$M = \text{term}$

$$\Sigma, L \times L$$



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



$N = \text{dokumen}$

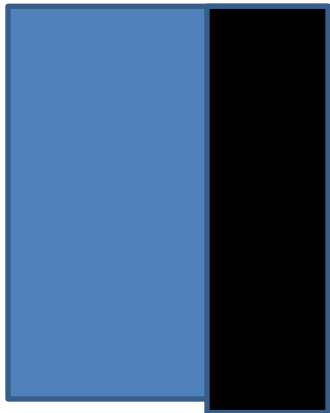
C_k = rank-k approximation of C

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

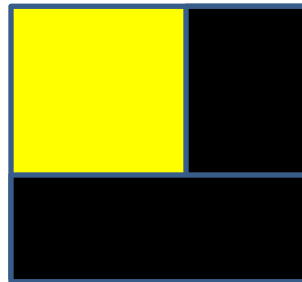
Jadikan NOL **L - K singular values paling kecil** di Σ !

basically dibuang L-K (kita mau ilangin Sigma kita jadi cuman K terbesar aja)

$$U_k, M \times K$$



$$\Sigma_k, K \times K$$



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



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

	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

← truncated

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

	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

Contoh: rank-2 approximation of C

C_2 adalah matriks **Rank-2** yang paling kecil error-nya dengan **C**.

LSA = SVD yang kita ambil K yang terbesar. Nah gunanya untuk reduksi dimensi

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

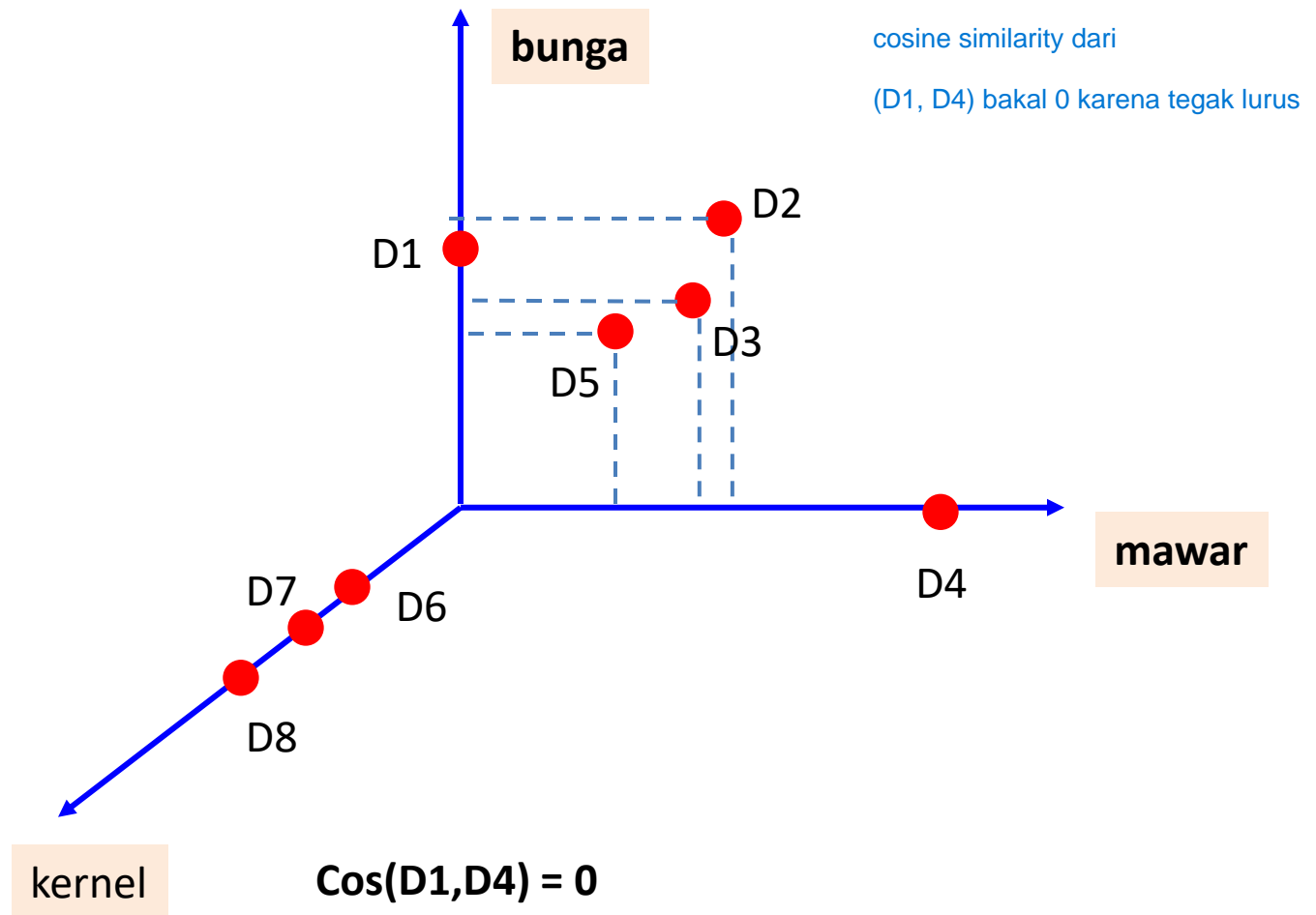
U_2

	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

Axis dengan eigenvalue terbesar pertama

kernel

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

D1

Axis dengan eigenvalue terbesar kedua

D2

D3

D5

dirotasi untuk variansi terbesar (atau fit dengan data-data yang ada)

D4

Di LSA kita buang yang eigen value paling jelek (eigen value paling kecil).

Nah artinya dibuang, dari 3D jadi 2D

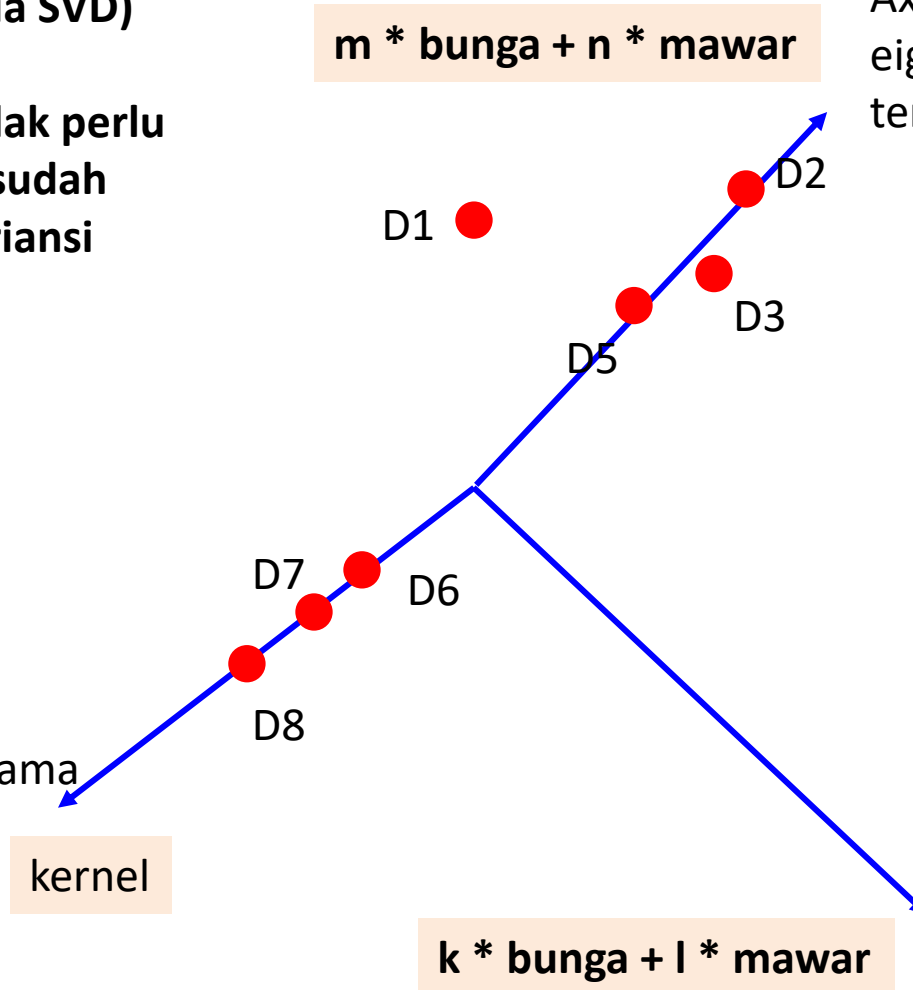
D7

D6

D8

$k * \text{bunga} + l * \text{mawar}$

Axis dengan eigenvalue terbesar ketiga



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

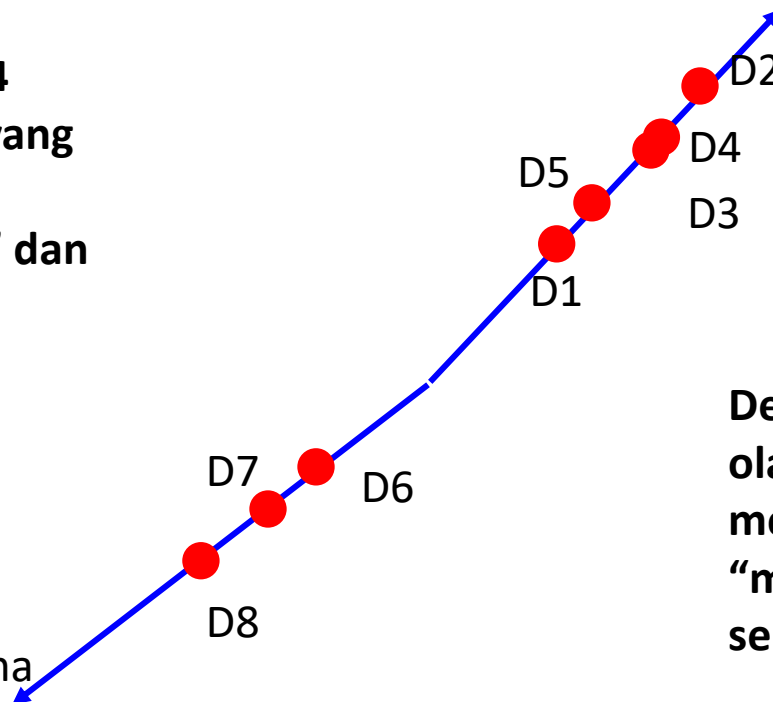
$m * \text{bunga} + n * \text{mawar}$

Axis dengan eigenvalue terbesar kedua

Nah tadi D1 dan D4 dirapetin karena dimensinya dikurangin. Oleh karena itu, $\text{Cos}(D1, D4)$ sekarang gak 0 lagi

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 $\Sigma_k \times V_k^T$

kalo di weighted

kalo unweighted ada di

V^T

Document Vector yang versi “unweighted”.

Dari Contoh sebelumnya:

2-Dimensional Document Embedding dari D3 adalah

Σ_2	18.93	0	\times	V_2^T	D3	$=$	D3
	0	14.49					
					0.30		4.35

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 $U_k \times \Sigma_k$

ini kalau weighted kalo unweighted ada di
U

Term Vector yang versi “unweighted”.

Dari Contoh sebelumnya:

2-Dimensional Document Embedding dari “Cancer” adalah

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

×

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:

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

Biasanya cosine similarity

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

Vektor kolom pada $\Sigma_k \times V_k^T$ yang terasosiasi dengan **d (jika ada)**; atau dokumen lain yang juga ditransformasi dengan cara yang sama.

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

Alternatif lain: seandainya kita mengasumsikan q adalah versi yang *unweighted* atau *normalized*:

$$q_k = \Sigma_k^{-1} \times U_k^T \times q$$

Biasanya cosine similarity

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

Vektor kolom pada V_k^T yang terasosiasi dengan **d (jika ada)**; atau dokumen lain yang juga ditransformasi dengan cara yang sama.

Contoh Soal LSA

Diberikan **Term-Document Matrix C** berikut dan hasil **SVD**-nya:

vektor dokumen ada di
 - weighted Sigma V^T
 - unweighted V^T

	D1	D2	D3
Bunga	2	1	0
Mawar	0	2	0
Kadal	0	0	3

C

	T1	T2	T3
Bunga	0	0.78	-0.61
Mawar	0	0.61	0.78
Kadal	1	0	0

U

	T1	T2	T3
T1	3	0	0
T2	0	2.56	0
T2	0		1.56

Σ

	D1	D2	D3
T1	0	0	1
T2	0.61	0.78	0
T3	-0.78	0.61	0

V^T

- Diberikan sebuah **Query = “mawar mawar”**; manakah yang relevansinya lebih tinggi antara **D1** dan **D3** pada **ruang vektor original**?
- Dengan Query yang sama manakah yang relevansinya lebih tinggi antara **D1** dan **D3** pada **ruang vektor “Rank-2” (i.e., Latent Semantic Indexing)**?
- Hitunglah cosine similarity antara Query sebelumnya dengan dokumen baru **D4 = {“Bunga”, “Kadal”}**!

Soal A) di ruang vektor original

$$Q = \begin{bmatrix} 0 \\ 2 \\ 0 \end{bmatrix} \quad D_1 = \begin{bmatrix} 2 \\ 0 \\ 0 \end{bmatrix} \quad D_3 = \begin{bmatrix} 0 \\ 0 \\ 3 \end{bmatrix}$$

$$\text{sim}(Q, D_1) = \frac{Q \cdot D_1}{\|Q\| \|D_1\|} = 0$$

$$\text{sim}(Q, D_3) = \frac{Q \cdot D_3}{\|Q\| \|D_3\|} = 0$$

D1 dan D3 sama-sama tidak mirip. **Make sense?**

Soal B) **Rank-2** Subspace

$$\Sigma_2 V_2^T = \begin{bmatrix} 0 & 0 & 3 \\ 1.57 & 2.01 & 0 \\ -1.23 & 0.96 & 0 \end{bmatrix}$$

dicoret



$$D'_1 = \begin{bmatrix} 0 \\ 1.57 \end{bmatrix}$$

$$D'_3 = \begin{bmatrix} 3 \\ 0 \end{bmatrix}$$

$$Q' = U_2^T Q = \begin{bmatrix} 0 \\ 1.23 \end{bmatrix}$$

Sebagai bahan renungan, apa yang terjadi jika masih tetap pada **Rank-3 Subspace**? Apakah **sim(Q, D1)** juga > 0? Mengapa?

$$\text{sim}(Q', D'_1) = \frac{Q' \cdot D'_1}{\|Q'\| \|D'_1\|} = \frac{1.93}{1.23 \times 1.57} = 1$$

$$\text{sim}(Q', D'_3) = \frac{Q' \cdot D'_3}{\|Q'\| \|D'_3\|} = 0$$

D1 mirip dengan Q; D3 tidak

Soal B) **Rank-2** Subspace (cara alternatif, **unweighted vector**)

$$V_2^T = \begin{bmatrix} 0 & 0 & 1 \\ 0.61 & 0.78 & 0 \\ -0.78 & 0.61 & 0 \end{bmatrix}$$

dicoret



$$D'_1 = \begin{bmatrix} 0 \\ 0.61 \end{bmatrix}$$

$$D'_3 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$Q' = \Sigma^{-1} U_2^T Q = \begin{bmatrix} 0 \\ 0.48 \end{bmatrix}$$

$$\text{sim}(Q', D'_1) = \frac{Q' \cdot D'_1}{\|Q'\| \|D'_1\|} = \frac{0.293}{0.48 \times 0.61} = 1$$

$$\text{sim}(Q', D'_3) = \frac{Q' \cdot D'_3}{\|Q'\| \|D'_3\|} = 0$$

D1 mirip dengan Q; D3 tidak

Soal C) Silakan kerjakan sendiri untuk Latihan (500 Point)

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 $M \times N$):

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

term topic matrix

- Informasi apa yang ada pada U ?
baris = Term
- Informasi apa yang ada pada Σ ? Apa makna singular values?
kolom = topic (basis baru yang terbuat/latent topic untuk DOKUMEN, peleburan topic tinggi dan topic berat)
matrix diagonal
- makna: seberapa penting topic yang ada
- Informasi apa yang ada pada V^T ?
baris = latent topic (basis baru utk TERM)
kolom = dokumen
- Bagaimana menghitung U , Σ , dan V^T ?
 $U = C C^T$
 $V = C^T C$
sigma = Akar dari eigen values
- Dimanakah *vector representation of words* berada?
term embedding di baris U
weighted U Σ V
- Bagaimana caranya mendapatkan *low-dimensional vector of words (dense)* berukuran $K (< \min(M, N))$?
 U ambil k docs dan potong sigma jadi $k \times k$
- Bagaimana caranya mendapatkan *low-dimensional vector of documents (dense)* berukuran $K (< \min(M, N))$?
mirip sama yang U tapi buat V

tujuan LSA: dokumen-dokumen yang mengandung kata-kata yang dekat, memiliki similarity yang tinggi.

LSA -> pakai SVD dan buang latent topic yang singular valuenya kecil. Kita sisakan aja singular value yang nilainya besar

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

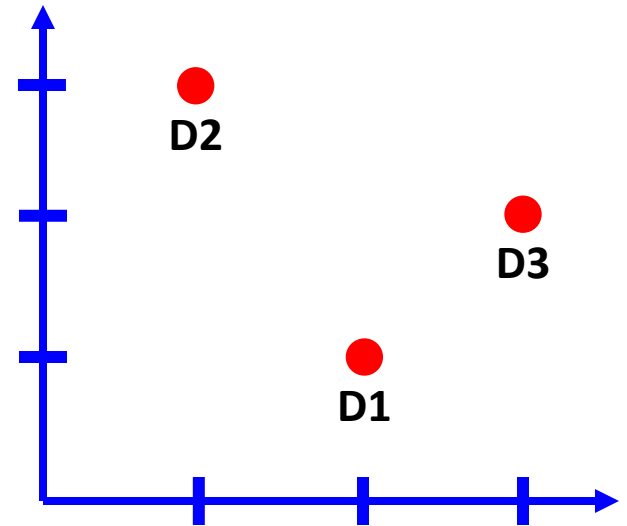
	D1	D2	D3	D4	D5
1	0.45	0.28	0.59	0.28	0.51
2	0.10	-0.59	0.30	-0.69	0.26
3	-0.11	-0.41	0.59	0.38	-0.56
4	-0.51	-0.42	-0.12	0.44	0.58
5	-0.70	0.46	0.42	-0.30	0.04

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

Simple Explanation of LSA

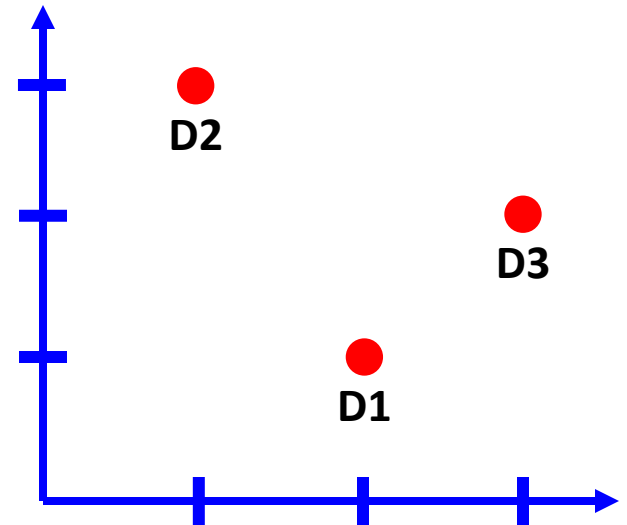
Misal, kita mempunyai term-document matrix berikut (bobot TF):

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



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

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



Misal, kita lakukan SVD:

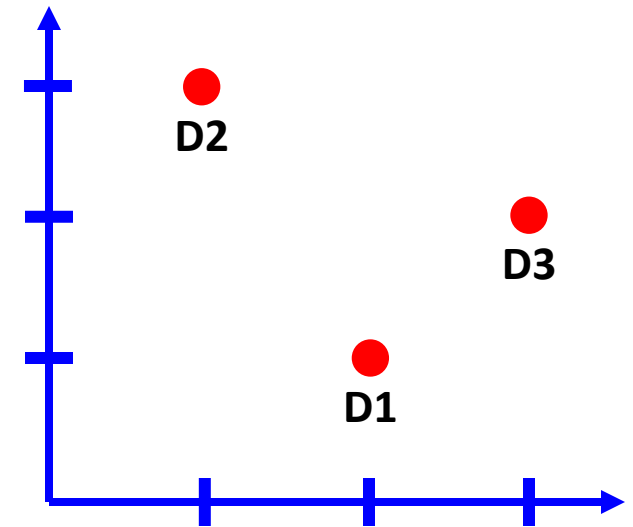
$$U = \begin{bmatrix} 0.71 & -0.71 \\ 0.71 & 0.71 \end{bmatrix}$$

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

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

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:

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

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

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

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

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):

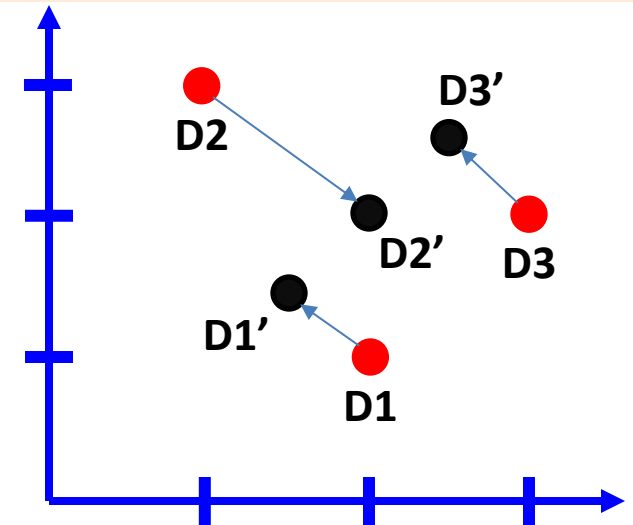
$$U = \begin{bmatrix} 0.71 & -0.71 \\ 0.71 & 0.71 \end{bmatrix}$$

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

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

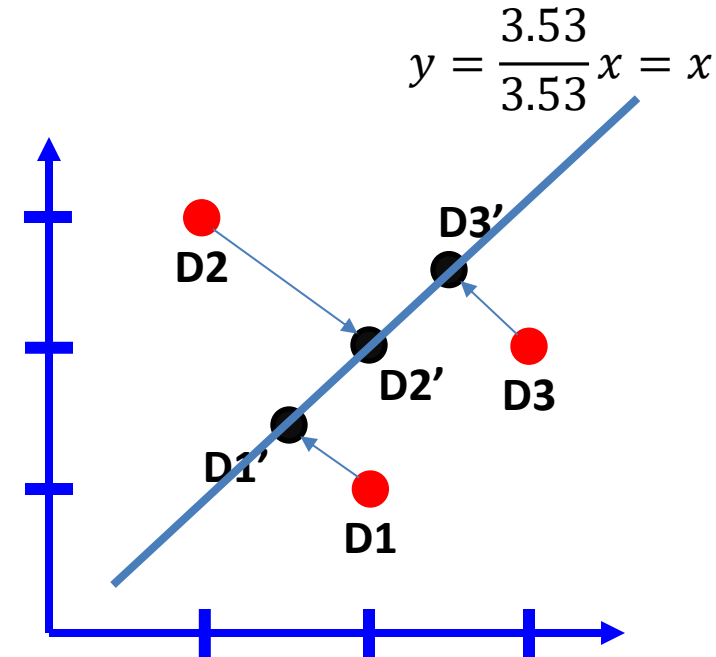
$$C_1 = U \times \Sigma_1 \times V^T = \begin{bmatrix} 1.5 & 2 & 2.5 \\ 1.5 & 2 & 2.5 \end{bmatrix}$$

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



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

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



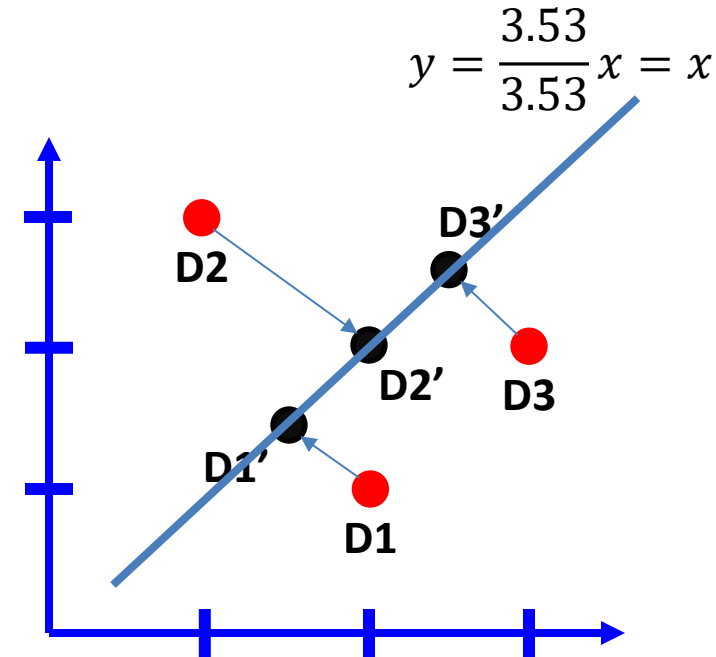
Vektor Kata:

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

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

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):

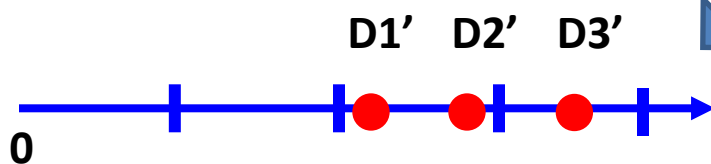
$$C = \begin{bmatrix} 2 & 1 & 3 \\ 1 & 3 & 2 \end{bmatrix}$$



Vektor Dokumen di Space Dim = 1:

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

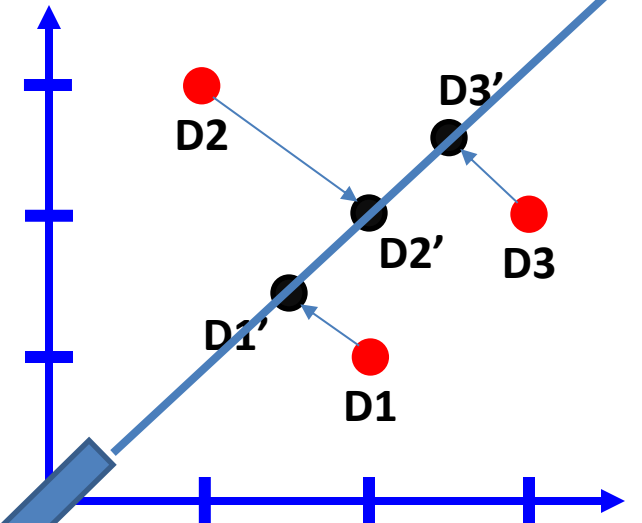
$$C_1 = U \times \Sigma_1 \times V^T = \begin{bmatrix} 1.5 & 2 & 2.5 \\ 1.5 & 2 & 2.5 \end{bmatrix}$$



Dimensi 1

Dimensi Original

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



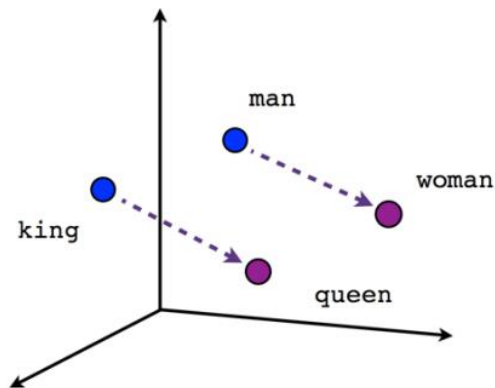
Neural Embeddings

Why Word Embeddings?

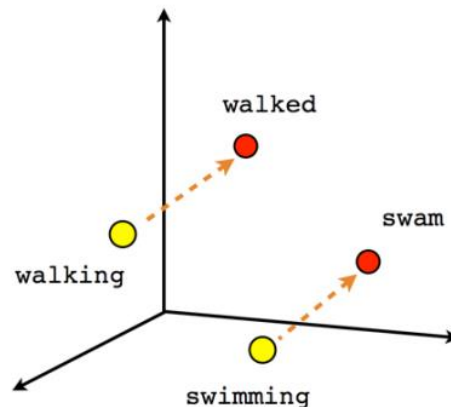
- Can capture the rich relational structure of the lexicon

kita bisa expand term term waktu user search query, misal kayak

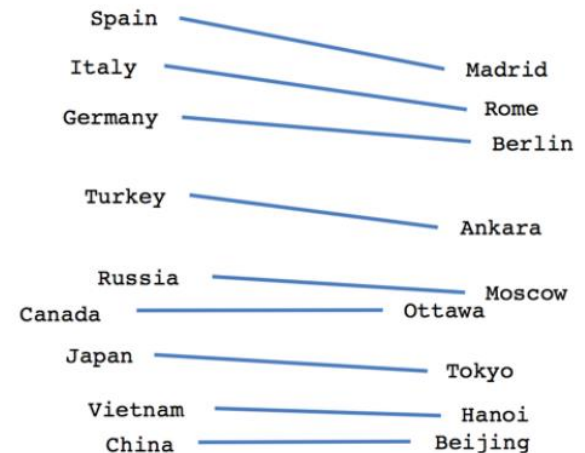
spain di expand madrid



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.

word embedding basis neural network

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**)

INPUT PROJECTION OUTPUT

- jadi lihat context dari sebelumnya aja (kata-kata sebelumnya)

-
- Diagram illustrating a Markov Decision Process (MDP) for word prediction. The states are represented by boxes, and transitions are indicated by arrows.
- Initial state: $w(t)$
 - Intermediate state (unlabeled box)
 - Transitions from the intermediate state lead to three possible next states: $w(t-2)$, $w(t-1)$, and $w(t+1)$.
 - Transitions from the intermediate state also lead to $w(t+2)$.
- Labels for states: $w(t)$, $w(t-2)$, $w(t-1)$, $w(t+1)$, $w(t+2)$.

- We seek a model for

$$P(w_{t+j} \mid w_t)$$

nah ini bisa pakai skip-gram

Skip-Gram

ini approach menggunakan NN (neural network)

Feed-Forward Process

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

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

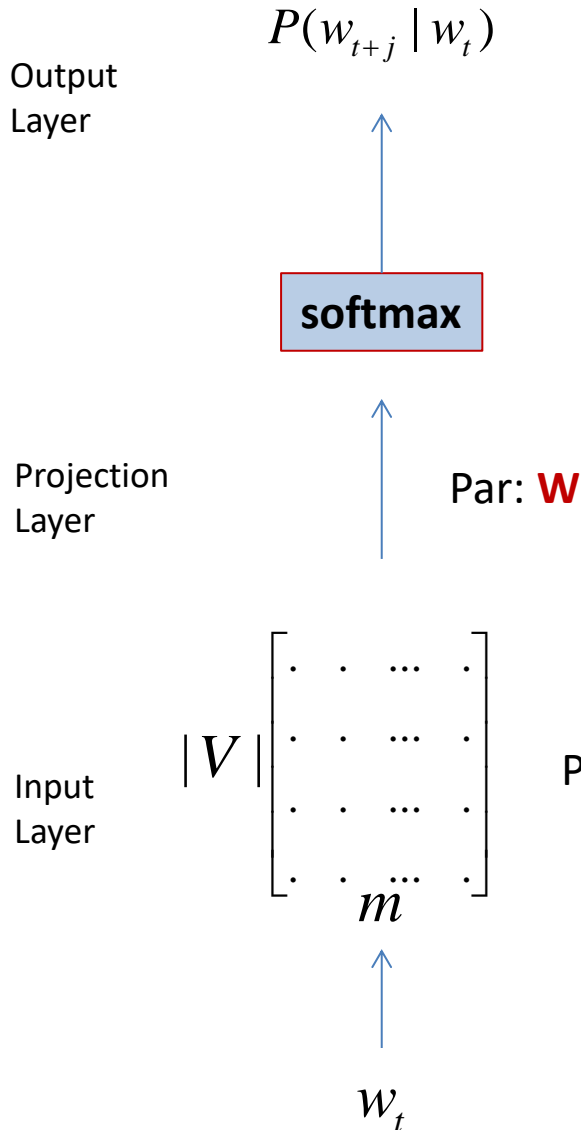
$$x = C(w_t)$$

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

word embedding tumor ada 2 pendekatan:

- ada yang pakai matrix center aja sebagai word embedding
- ada juga pendekatan lain, yaitu rata-rata embedding C dan rata-rata di matrix W (context word)

Sudut Pandang Lain

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

|Vocab| X |Embedding Dimension|

Matriks C

Center word

embeddings of tumor

**Low-dimensional Dense
Embeddings of Tumor**

setiap baris maksudnya
word embedding

$$\begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}^T \times \begin{bmatrix} 0.1 & 0.2 & 0.2 \\ 0.1 & 0.1 & 0.3 \\ 0.8 & 0.4 & 0.1 \\ 0.9 & 0.3 & 0.1 \end{bmatrix} = \begin{bmatrix} 0.1 \\ 0.1 \\ 0.3 \end{bmatrix}^T$$

berapa banyak kata di vocab
x
dimensi embedding
(hyperparametercoba-coba)

One-hot
representation of
tumor

yang di blue highlight:

adalah embedding suatu kata sebagai CENTER

$$P(\textit{rose}|\textit{tumor}) = 0.25$$

|Embedding Dimension| X |VOCAB|

Matriks W

Context word

embeddings of rose

$$\textit{softmax} \left\{ \begin{bmatrix} 0.1 \\ 0.1 \\ 0.3 \end{bmatrix}^T \times \begin{bmatrix} 0.3 & 0.2 & 0.3 & 0.1 \\ 0.2 & 0.2 & 0.4 & 0.7 \\ 0.2 & 0.3 & 0.3 & 0.1 \end{bmatrix} \right\} = \begin{bmatrix} 0.24 \\ 0.25 \\ 0.26 \\ 0.25 \end{bmatrix}$$

Skip-Gram

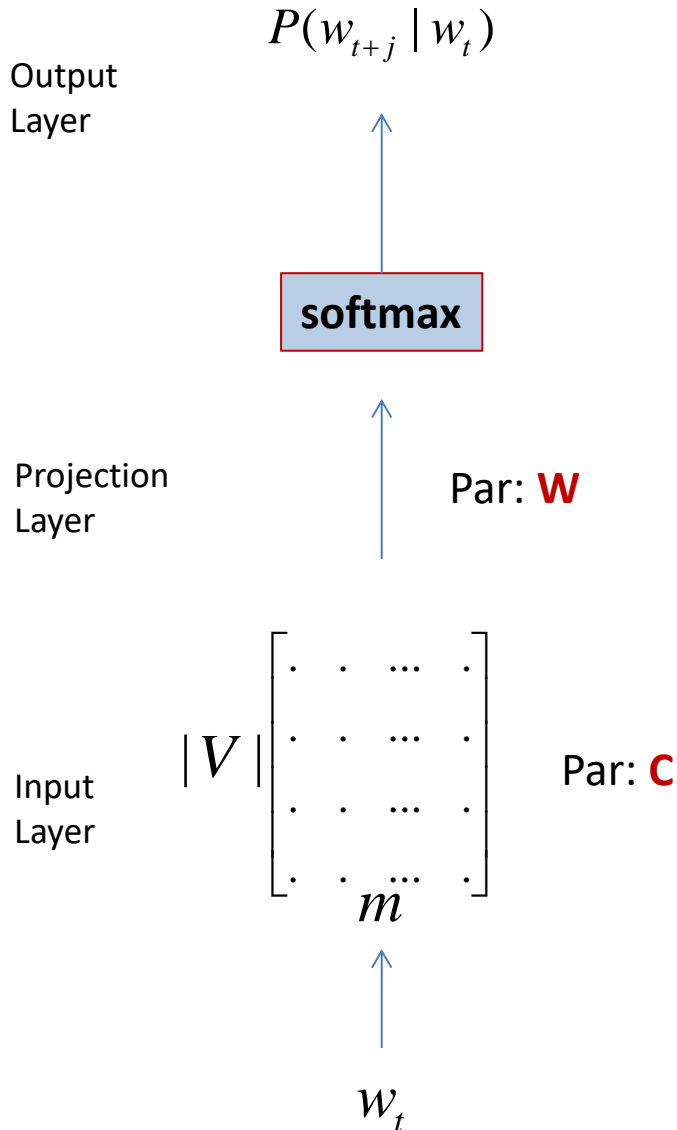
Dengan kata lain, Skip-Gram juga bisa dinyatakan dengan:

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

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

Baris di matriks C

Kolom di matriks W

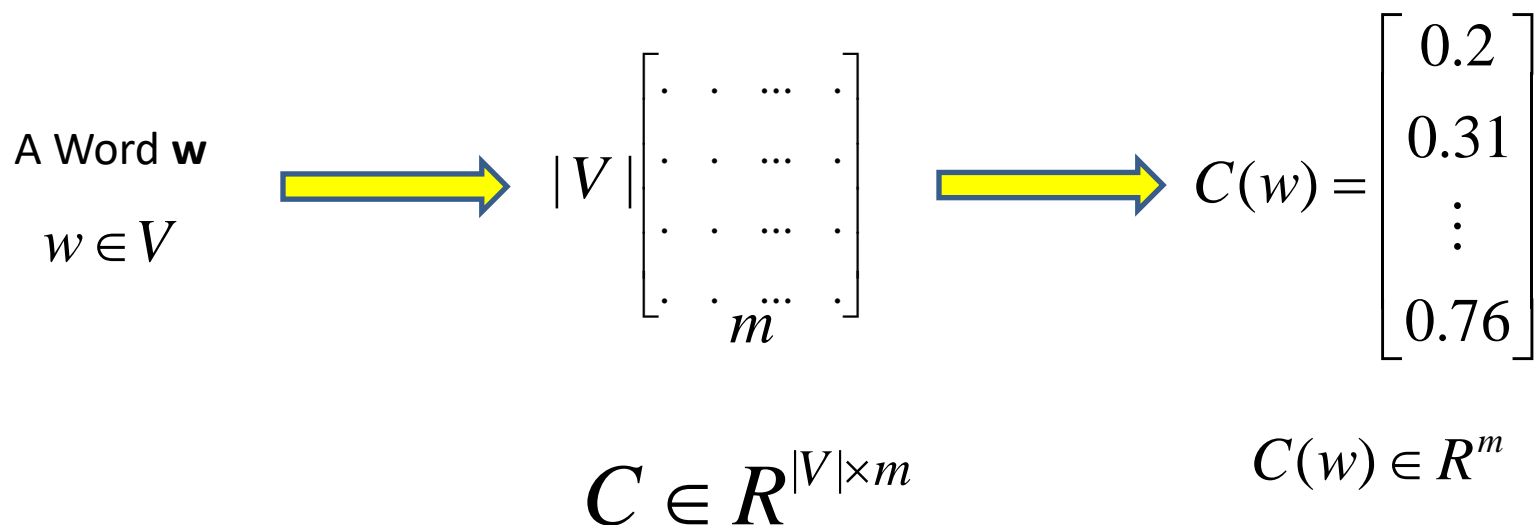


Where are the Word Embeddings?

- 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:

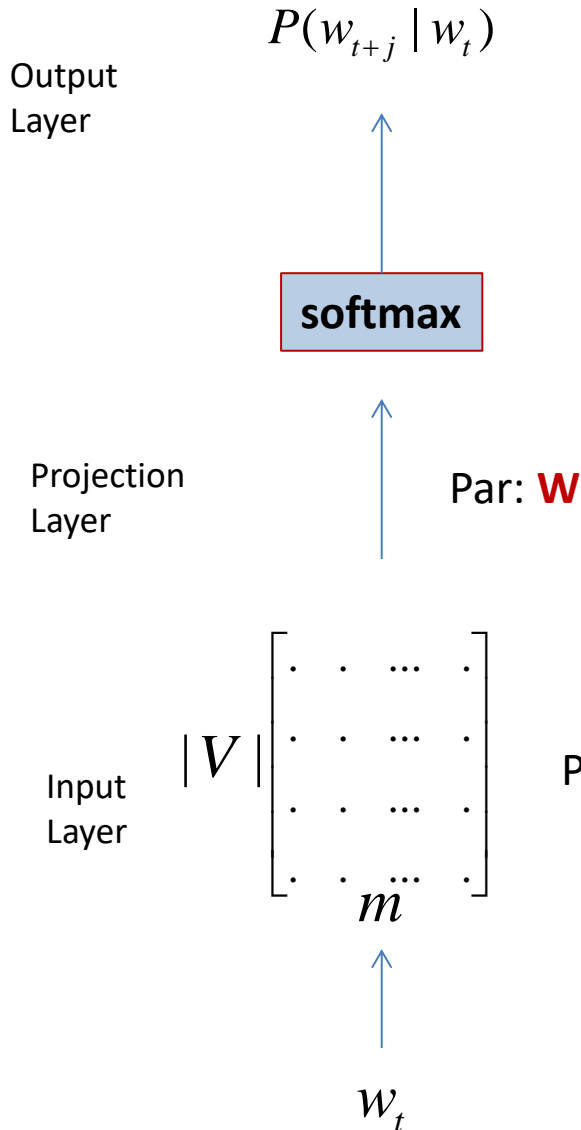
Given training data $w_1, w_2, w_3, \dots, w_{T-1}, w_T$

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

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

Regularization Terms

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



Skip-Gram

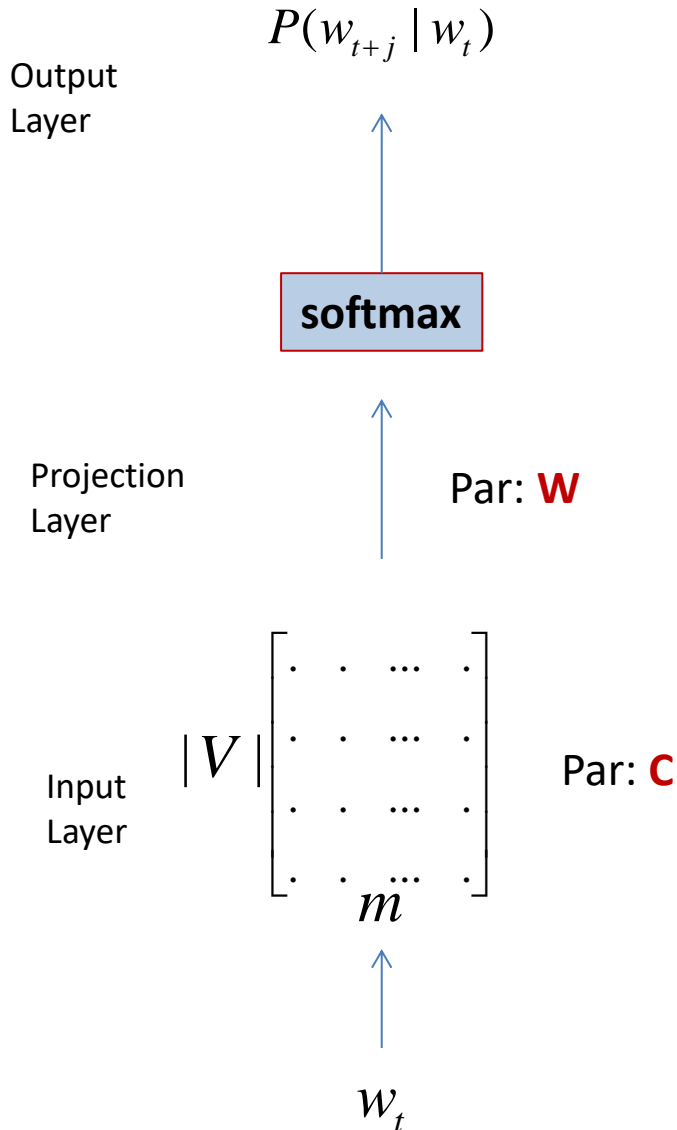
Training

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.

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

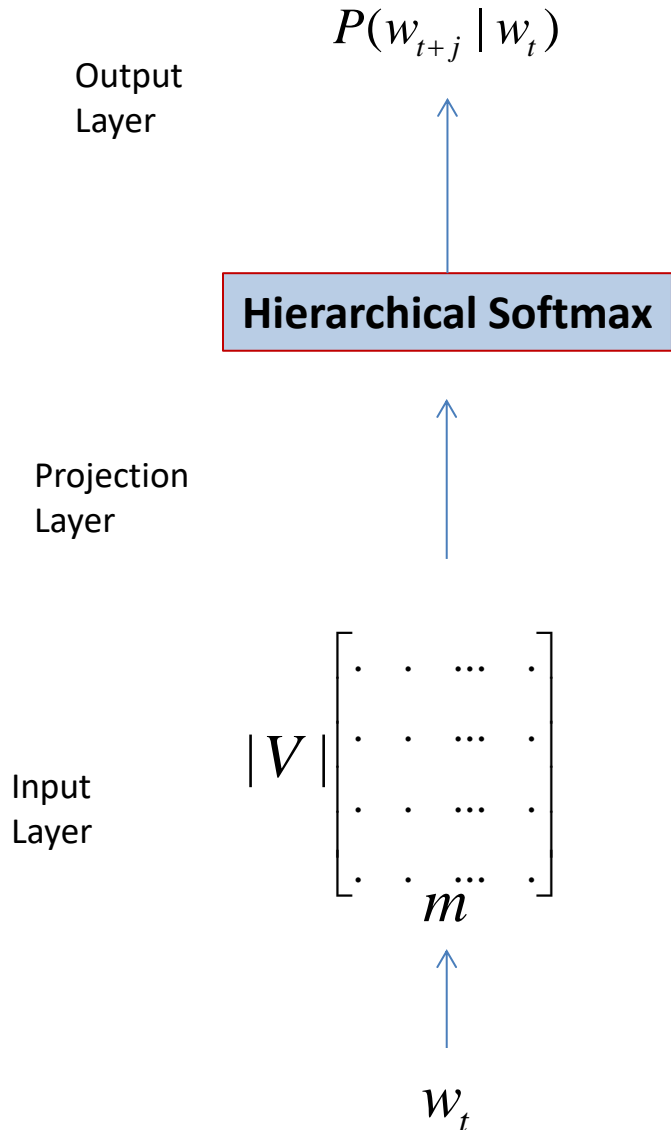
D is the maximum distance of the words

biar komputasi softmax di ujung bisa lebih murah

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

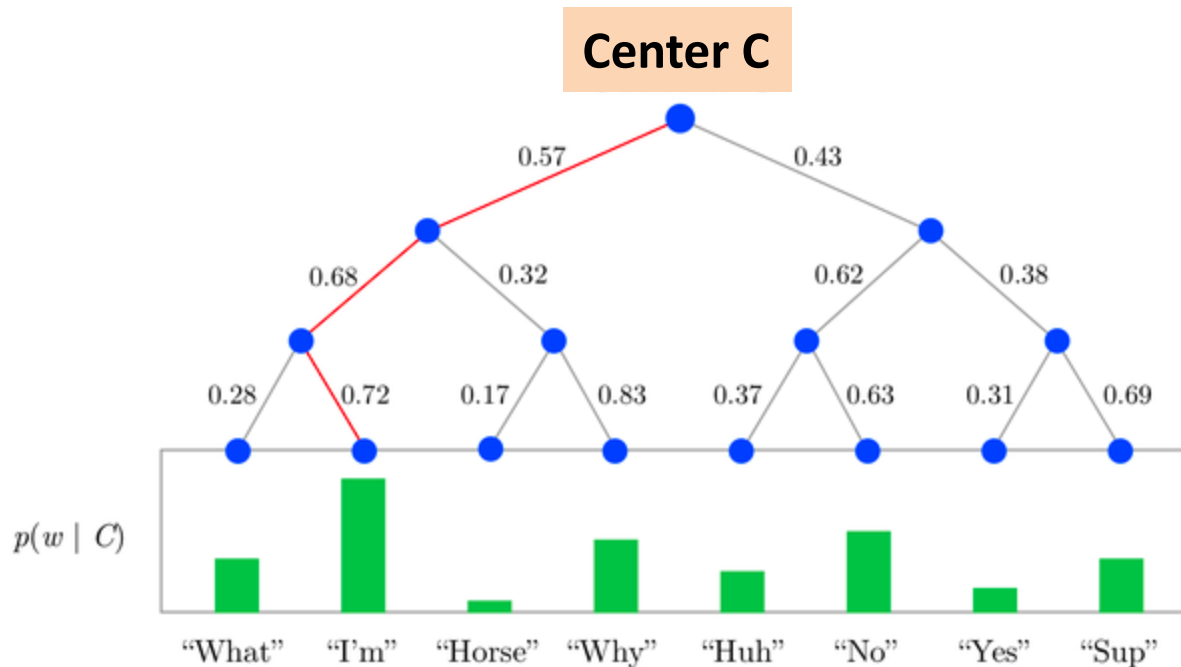
$$Q = D \times (m + m \times \log_2 |V|)$$

(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

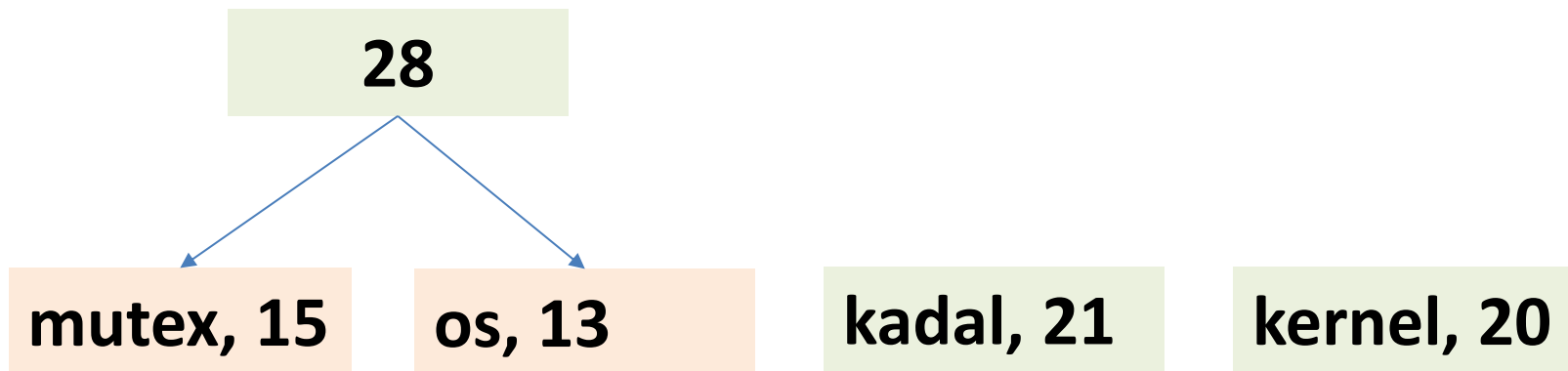
os, 13

kernel, 20

Hierarchical Softmax

Huffman Tree

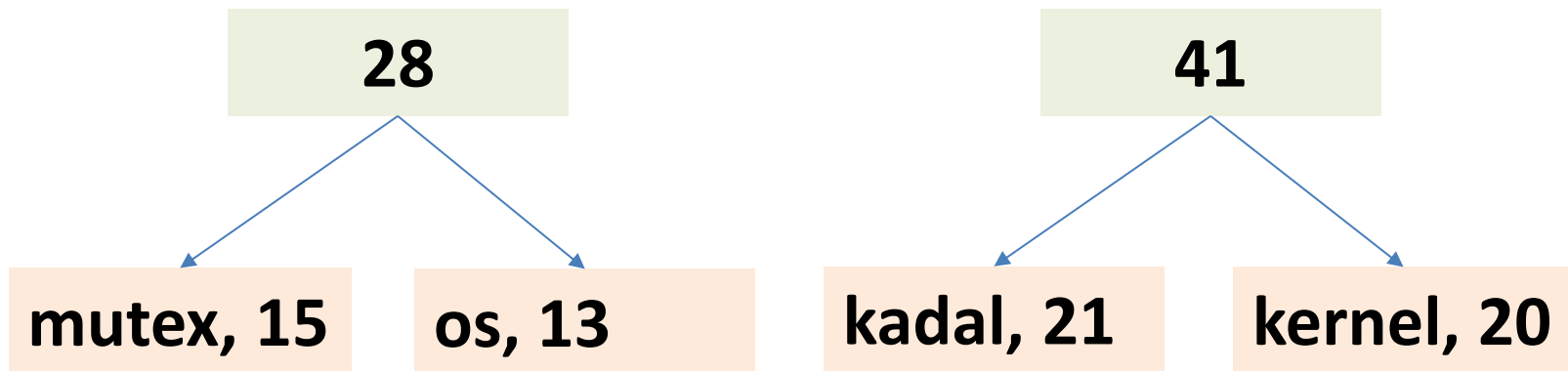
Gabung dua buah kata dengan frekuensi paling kecil; dan proses ini dilakukan terus.



Hierarchical Softmax

Huffman Tree

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

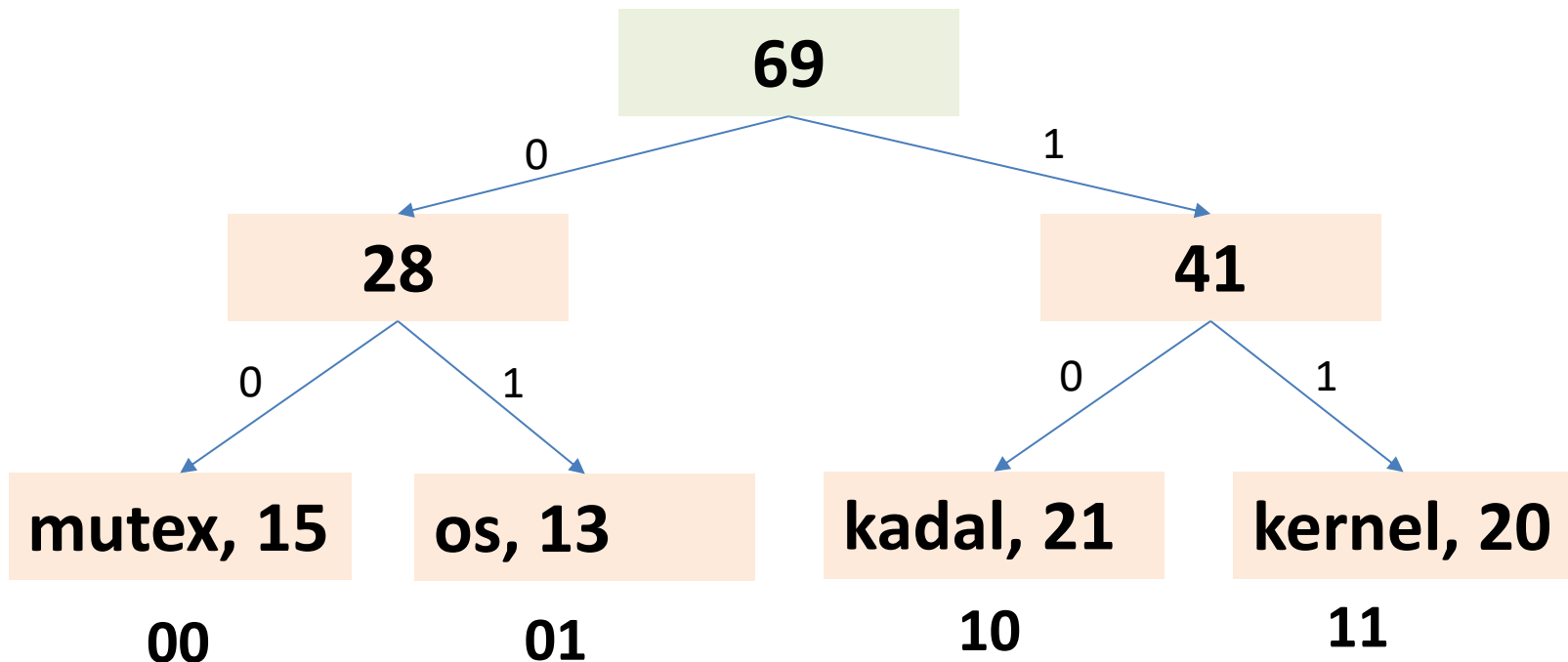


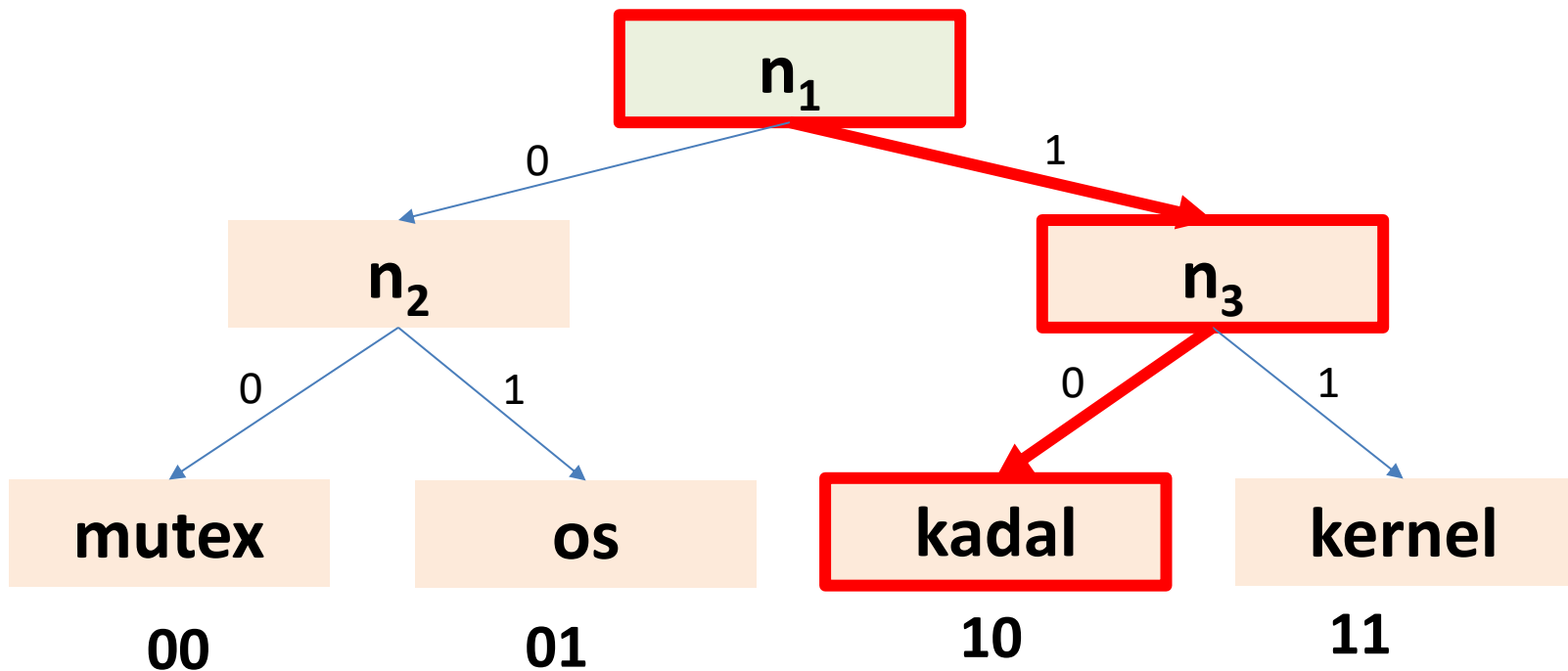
Hierarchical Softmax

Huffman Tree

Assign kode biner ke setiap kata di vocabulary

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





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

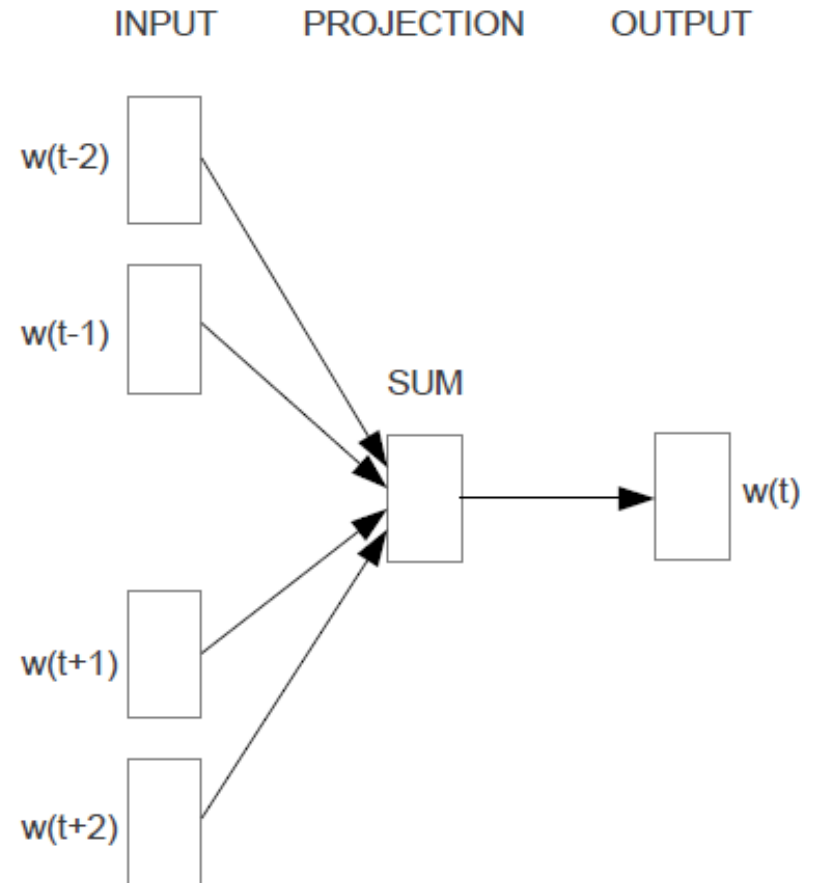
$$P_n(\text{left}|c) = \sigma(v_n^T \cdot \text{center}(c))$$

Trainable vector for node n

Hidden layer representation of the center c

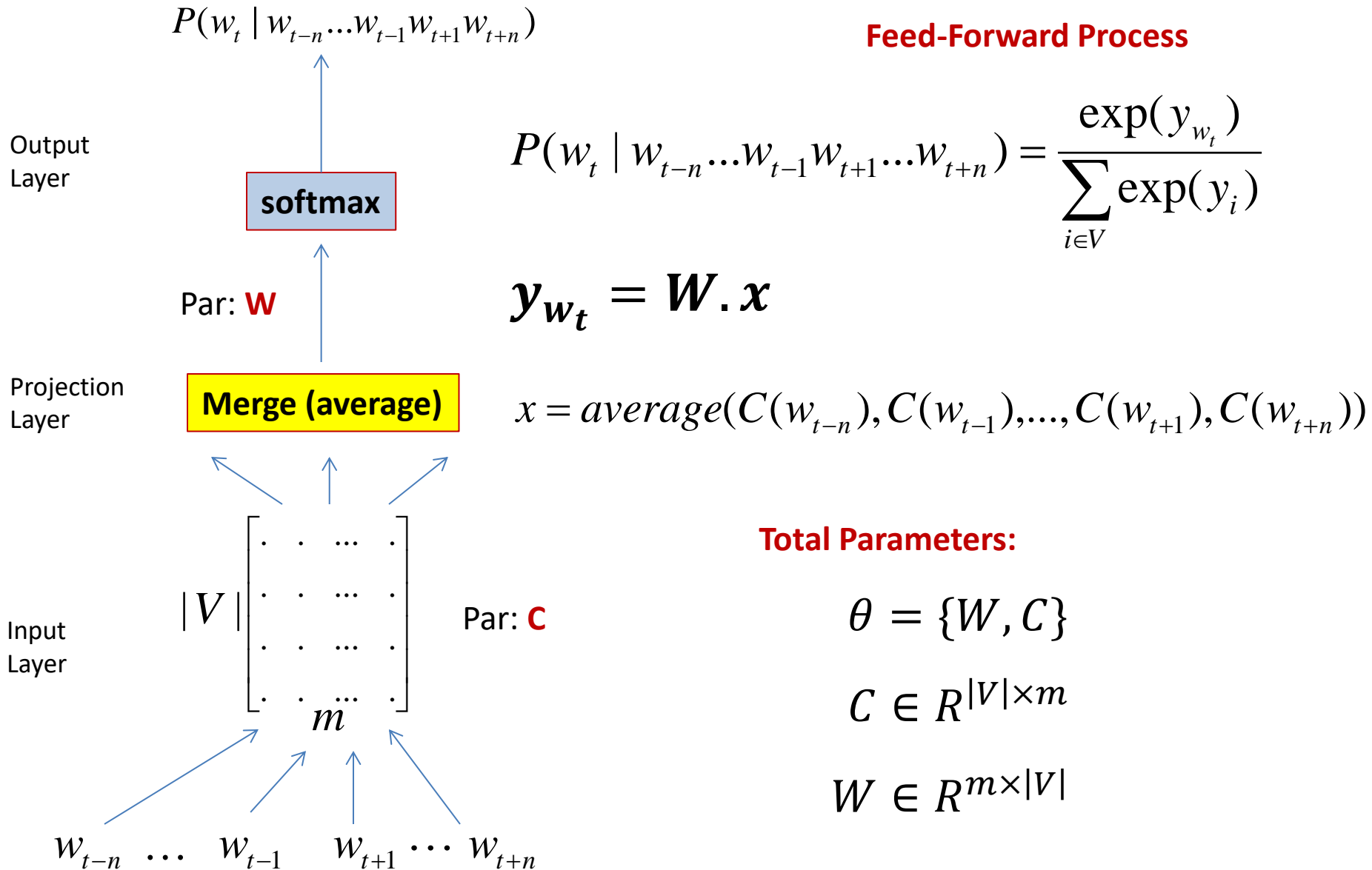
Continuous Bag-of-Words

- 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**.



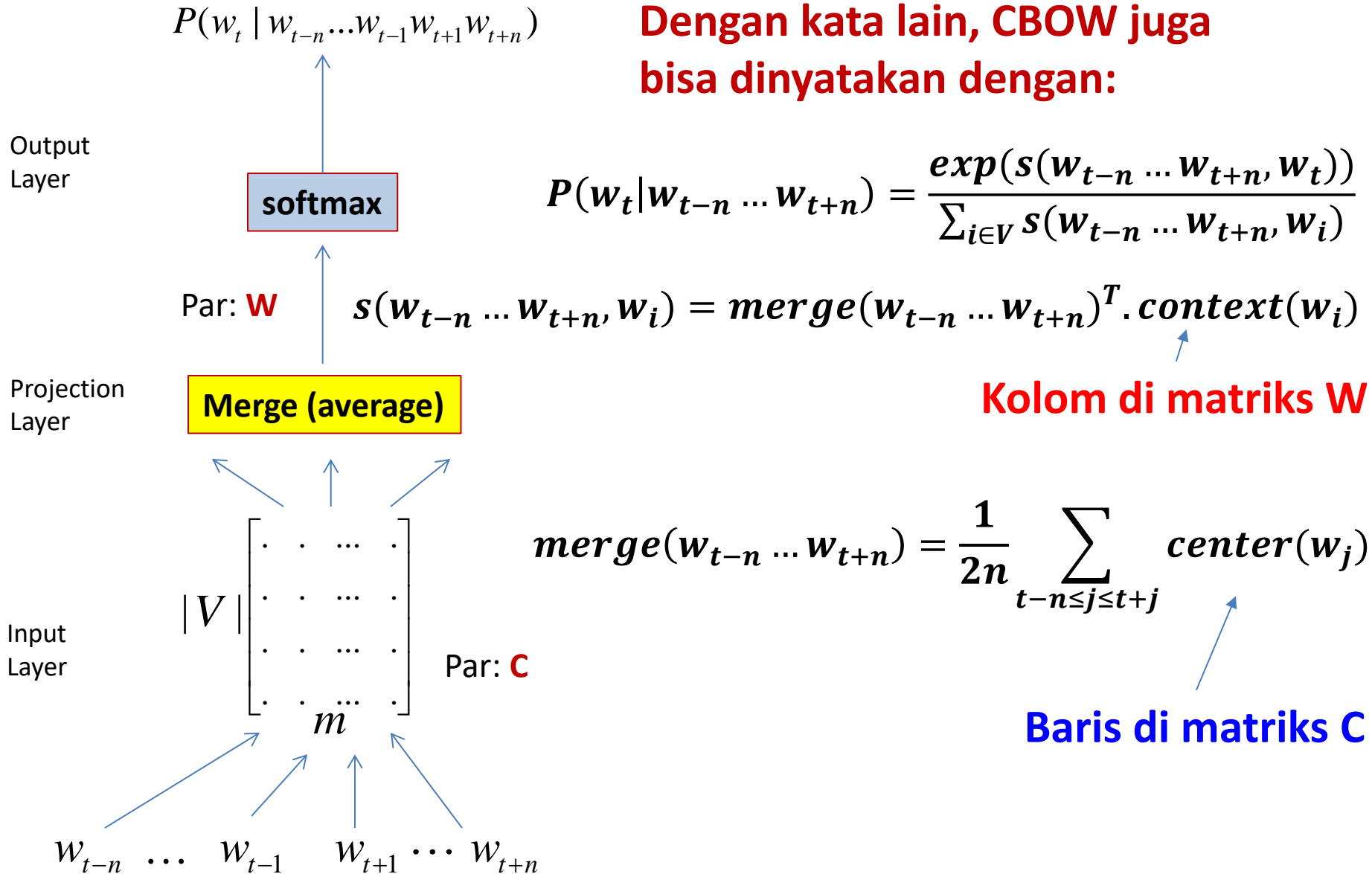
We seek a model for $P(w_t \mid w_{t-n} \dots w_{t-1} w_{t+1} \dots w_{t+n})$

Continuous Bag-of-Words

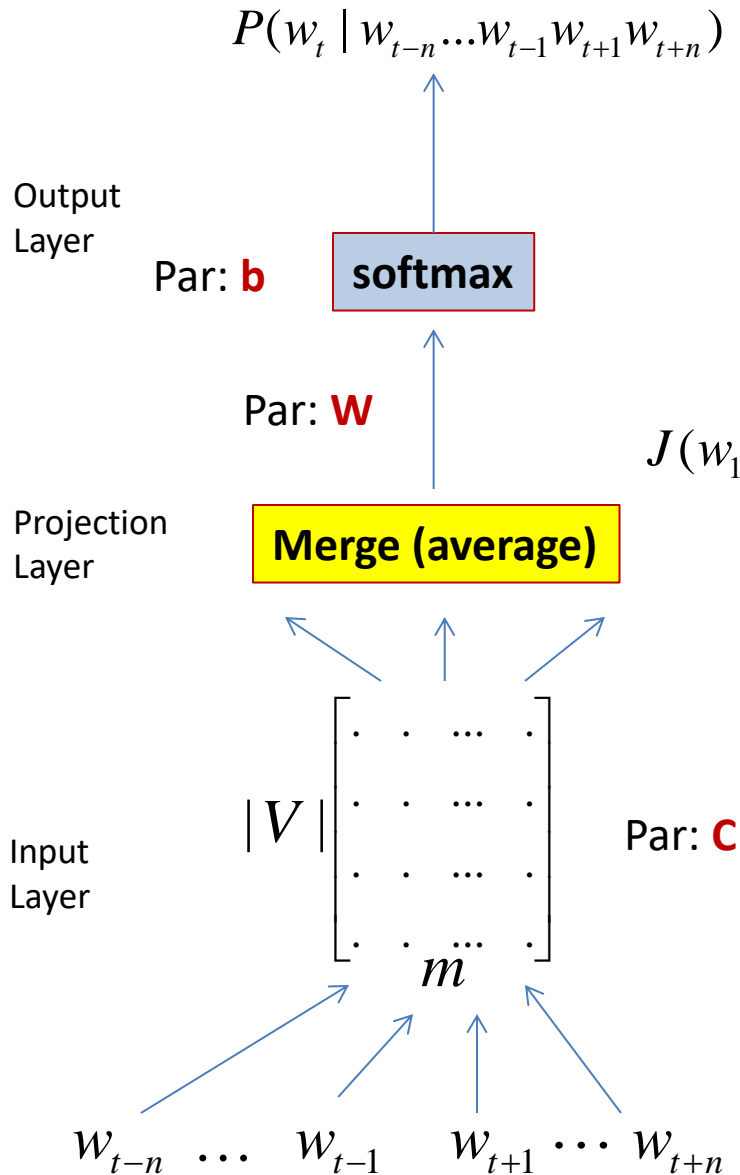


Continuous Bag-of-Words

Dengan kata lain, CBOW juga bisa dinyatakan dengan:



Continuous Bag-of-Words



Training

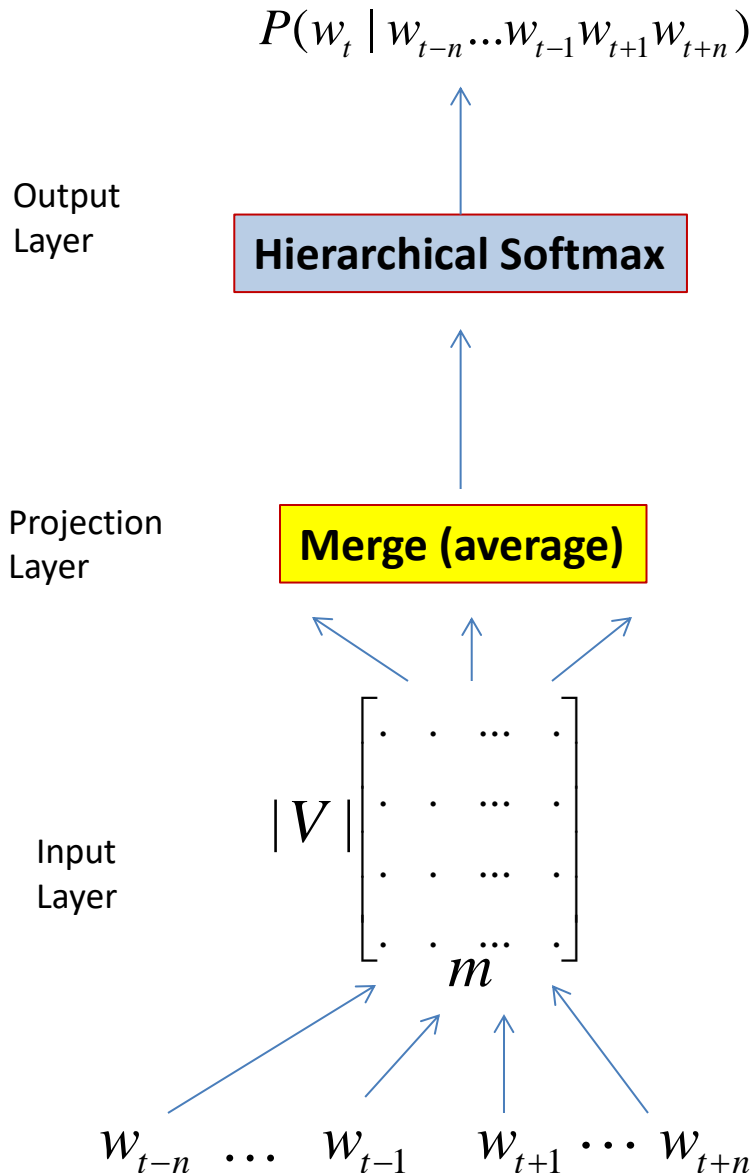
Training is achieved by looking θ that maximizes the following Cost Function:

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



Training

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

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

biar komputasi softmax di ujung bisa lebih murah

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 -> $z_{g1} = [0.3, 0.1, \dots]$

– whe -> $z_{g2} = [0.1, 0.4, \dots]$

– her -> $z_{g3} = [0.5, 0.9, \dots]$

– ere -> $z_{g4} = [0.5, 0.5, \dots]$

– re> -> $z_{g5} = [0.3, 0.1, \dots]$

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

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