

DISTRIBUTIONAL SEMANTICS AND WORD2VEC

L5-1

- BOW Feature Representation is pretty ONE OF CONVENIENCE! IT GIVES US A WAY TO REPRESENT VARIABLE LENGTH SEQUENCES ~~DATA~~ WITH A FIXED INPUT SIZE.
- UNFORTUNATELY, LANGUAGE IS CONTEXTUAL. THE MEANING OF A WORD IS A FUNCTION OF ITS ~~THE~~ CONTEXT.
- ANY MODEL of NATURAL LANGUAGE THEREFORE MUST ACCOUNT FOR THE MEANING OF WORDS IN CONTEXT.
 \therefore we must move away from BOW!
- NEW FEATURE REPRESENTATION:

VOCABULARY: $X = I_N = N \begin{bmatrix} 1 & \dots & 0 \\ 0 & \dots & 1 \end{bmatrix}$ $N = \text{VOCAB SIZE}$

EXAMPLE

$\therefore X_i = [0 \ 0 \ 0 \ 1 \ 0] = \text{ONE-HOT VECTOR REPRESENTATION of } i^{\text{th}} \text{ word in } X.$

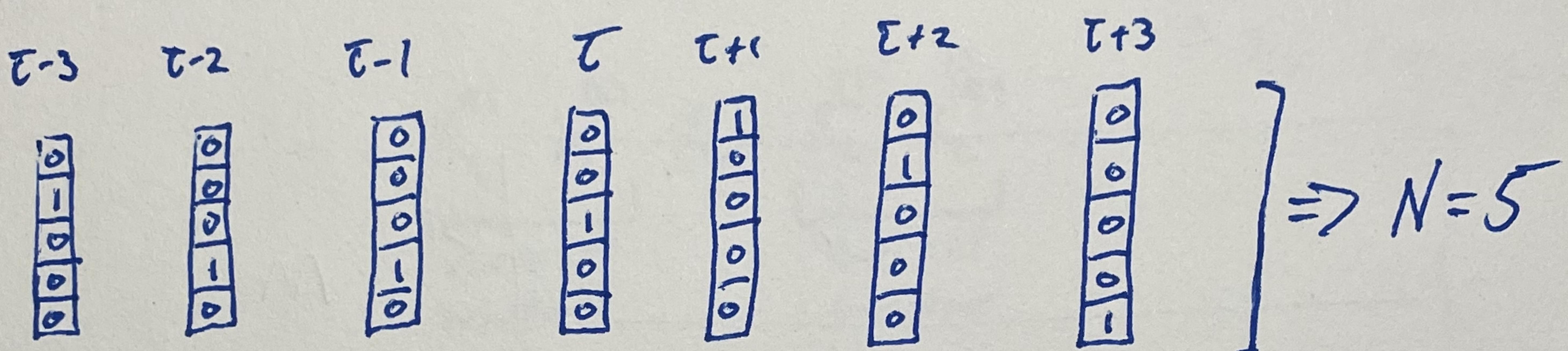
SEQUENCE-PRESERVING FEATURE REPR.

LS-2

↳ NEW FEATURE REPR. :

$\dots X^{(t-3)} \ X^{(t-2)} \ X^{(t-1)} \ X^{(t)} \ X^{(t+1)} \ X^{(t+2)} \ X^{(t+3)} \ \dots \dots \dots$

↳ t INDEXES THE SEQUENCE DIMENSION



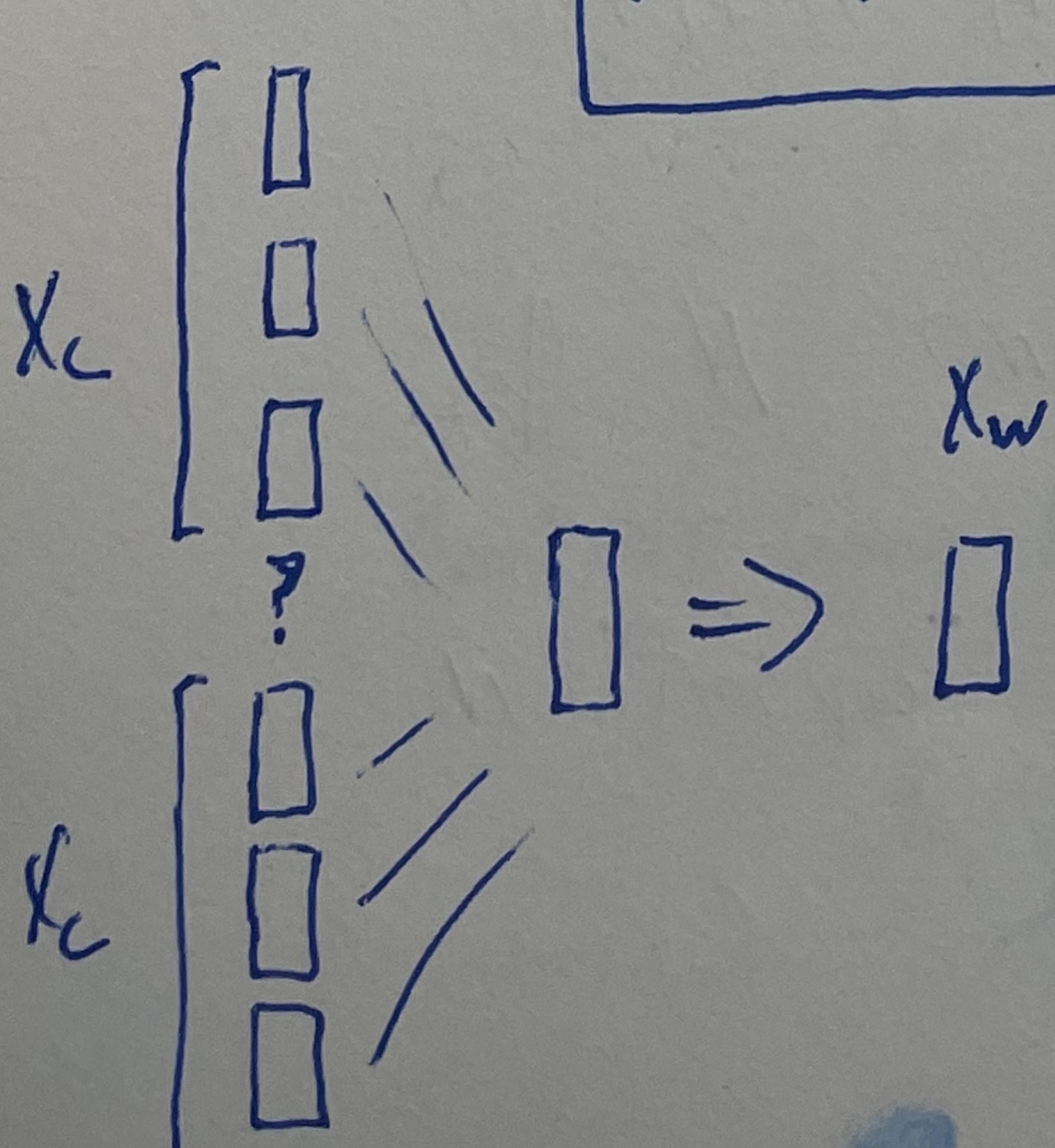
↳ SEQUENCE OF ONE-HOT ENCODED WORDS

LANGUAGE MODELING APPROACHES

x_w = CENTER WORD
 x_c = CONTEXT WORD

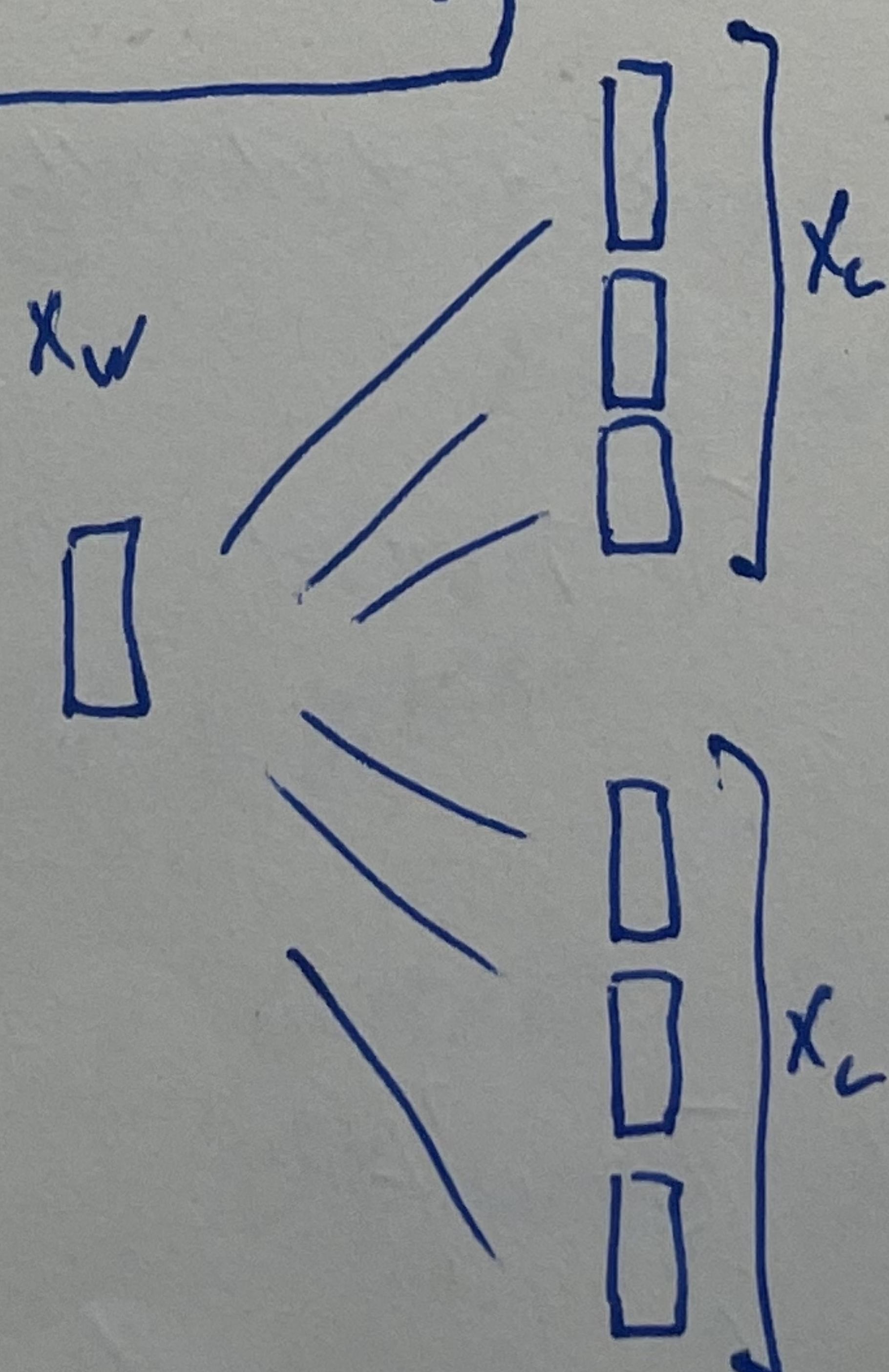
CONTINUOUS BAG OF WORDS

CBoW
 $P(x_w | x_c^{(1)}, \dots, x_c^{(c)})$



Skip-Gram

$P(x_c^{(i)} | x_w)$



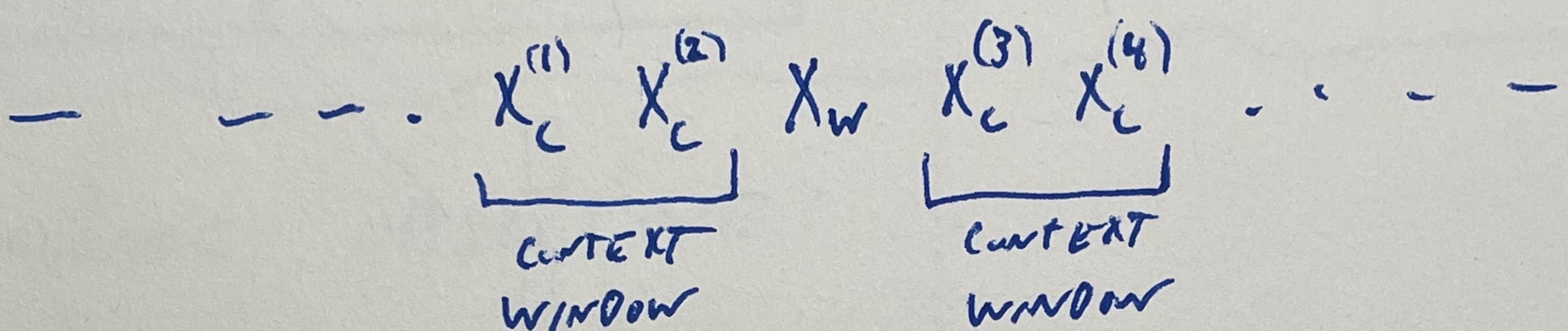
SKIP GRAM WORD2vec

L5-3

DATASET: $D = \{ X_w^{(1)}, X_c^{(1)} \dots X_w^{(M)}, X_c^{(M)} \}$

↳ M does NOT EQUAL number of words in CORPUS!

↳ For EACH (CENTER) word in CORPUS:



↳ $(X_w, X_c^{(1)}), (X_w, X_c^{(2)}), (X_w, X_c^{(3)}), (X_w, X_c^{(4)})$

$$M \approx \text{Num words} \times \text{context window} \times 2$$

Word2vec

CENTER word EMBEDDINGS: $U \in \mathbb{R}^{K \times N}$ $K \ll N$

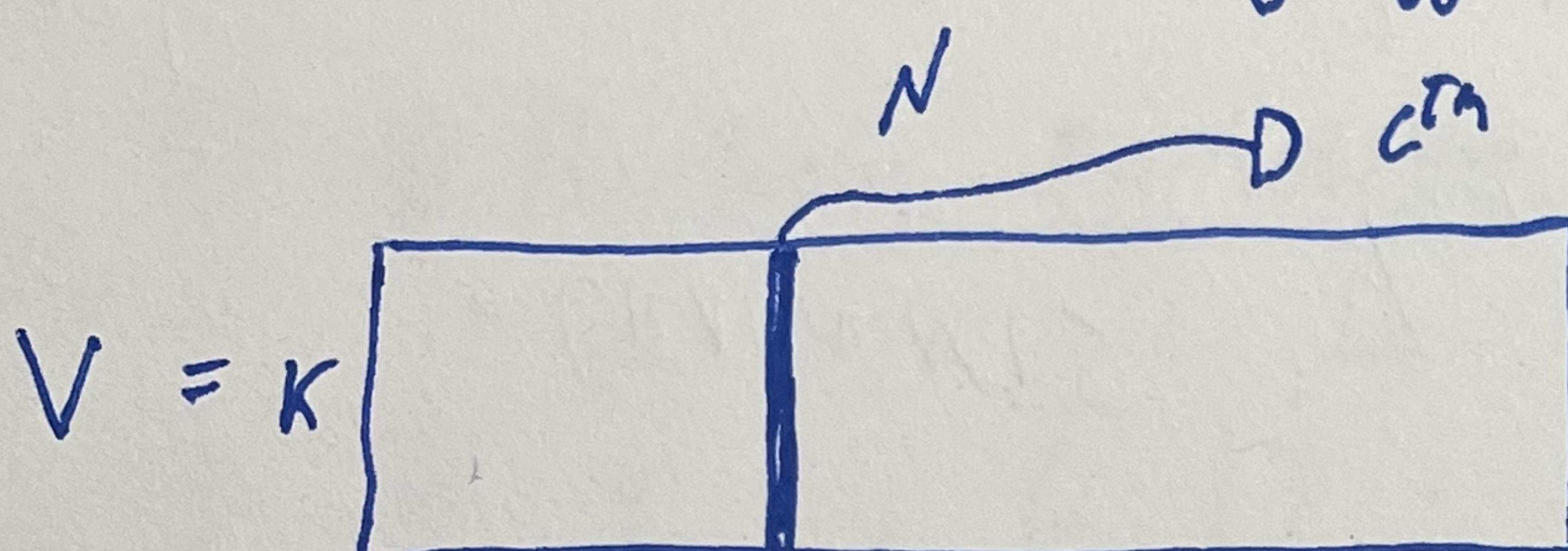
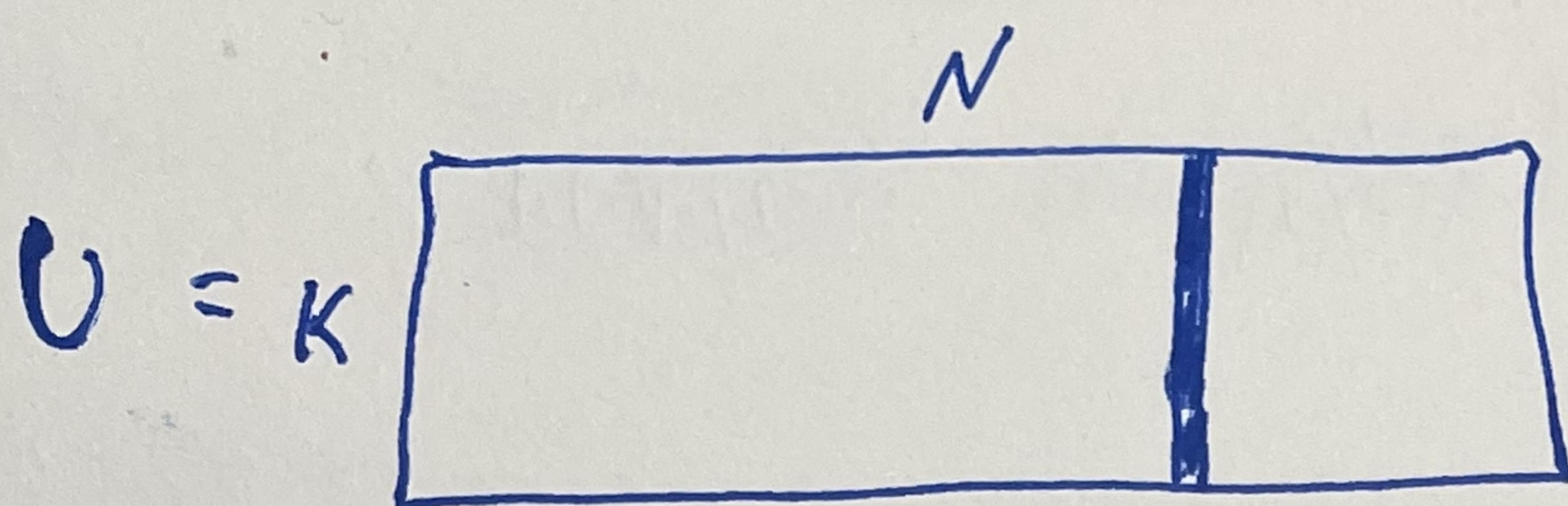
↳ $X_w \cdot U = u_w \in \mathbb{R}^{K \times 1}$ = CENTER word EMB
FOR w^{th} word in X

CONTEXT word EMBEDDINGS: $V \in \mathbb{R}^{K \times N}$

$X_c \cdot V = v_c \in \mathbb{R}^{K \times 1}$ = CONTEXT word EMB for
 c^{th} word in X .

Skip-Gram Word2vec Cont...

LS-4



$\hookrightarrow w^m$ column of $U = x_w \cdot U = u_w$

c^m column of $V = x_c \cdot V = v_c$

\hookrightarrow Word2vec OBJECTIVE:

- ① MAKE $u_w^T v_c$ LARGE if x_c IS IN CONTEXT of x_w IN D
- ② MAKE $u_w^T v_c$ SMALL if x_c IS NOT IN CONTEXT of x_w IN D

\hookrightarrow How Does Word2vec Do This?

(1) ANSWER: MLE of THE DISCRIMINATIVE MODEL $P(x_c | x_w)$

$$\hookrightarrow P(x_c | x_w; U, V) = \frac{e^{u_w^T v_c}}{\sum_{j=1}^N e^{u_w^T v_j}} \in [0, 1]^N$$

$P_{x_c | x_w}$ for Short

SKIP-GRAM Word2Vec ALGORITHM

LS-5

Goal: Minimize $NLL(U, V | D)$ using stochastic GD

ALGORITHM:

① Randomly Initialize $U, V \in \mathbb{R}^{k \times N}$

② $\forall x_c, x_w \in D:$

- "Lookup" the embedding u_w in U pointed to by x_w .

• COMPUTE: $z = u_w^T V \in \mathbb{R}^N$

• COMPUTE: $P_{x_c|x_w} = \frac{e^z}{\sum_{j=1}^N e^{z_j}}$ ~~normalize~~ $\in \{0, 1\}^N$

• COMPUTE: $NLL = -x_c \cdot \log P_{x_c|x_w} \in \mathbb{R}$

• COMPUTE: $\nabla_{u_w} NLL = V \cdot (P_{x_c|x_w} - x_c)^T \in \mathbb{R}^{k \times 1}$

• COMPUTE: $\nabla_v NLL = u_w^T \cdot (P_{x_c|x_w} - x_c) \in \mathbb{R}^{k \times N}$

• UPDATE: $(U^T)_w = (U^T)_w - \eta \nabla_{u_w} NLL$

• UPDATE: $V = V - \eta \nabla_v NLL$

REPEAT