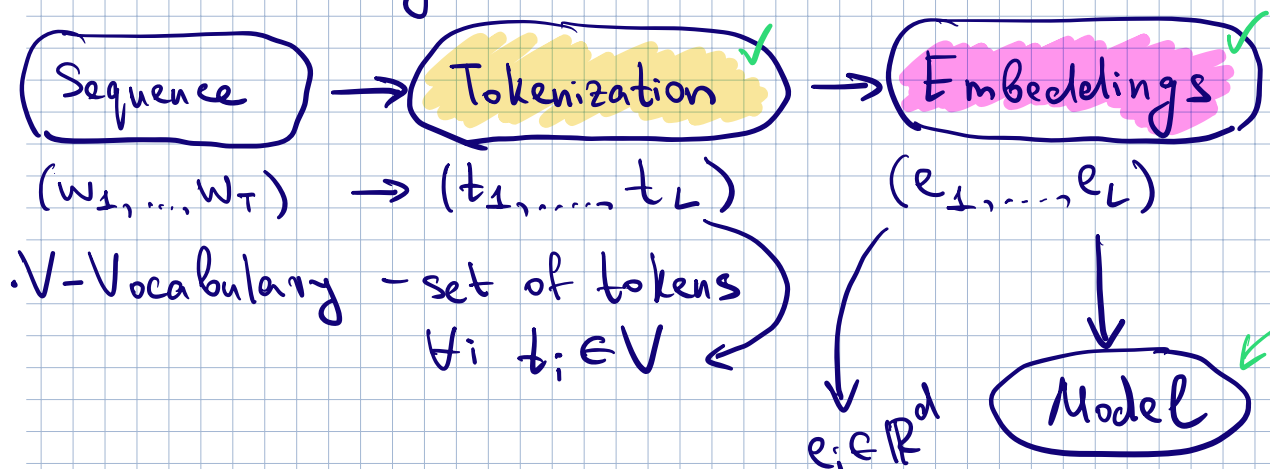


$\begin{matrix} b & o & s & p & t \\ s & o & s & p & t \end{matrix}$
 $\begin{matrix} g & o & s & p & t \\ n & o & n & n & t \end{matrix}$
,
 $\begin{matrix} g & o & s & p & o & t & o & \bar{u} \end{matrix}$

Sequence processing

Processing pipeline

- Texts are discrete
- Variable length



Tokenization

• Character-level

$\begin{matrix} b & o & s & p & t \\ s & o & s & p & t \end{matrix}$
 $\begin{matrix} g & o & s & p & t \\ n & o & n & n & t \end{matrix}$
,
 $\begin{matrix} g & o & s & p & o & t & o & \bar{u} \end{matrix}$

- + Very small vocab
- + Easy to encode new words
- Lose language structure
- Too long sequences

• Word-level

|b|o|s|p|t| |g|o|s|p|t|,| |g|o|s|p|o|t|o|u| |
|s|o|s|p|t| |n|o|l|h|h|.|

- Very large vocab

- No tokens for new words

• <UNK> - token for unknown words

+ Preserve language structure

+ Shorter sequences

• Subword-tokenization

Byte-pair encoding (BPE)

|b|o|s|p|t| |g|o|s|p|t|,| |g|o|s|p|o|t|o|u| |
|s|o|s|p|t| |n|o|l|h|h|.|

1) Start with character-level tokenization

2) While $\text{size}(V) < \text{vocab_size}$

• Find most frequent pair of tokens

• Merge 2 tokens, add to vocab

|b|o|s|p|t| |g|o|s|p|t|,| |g|o|s|p|o|t|o|u| |
|s|o|s|p|t| |n|o|l|h|h|.|

0) {s, o, p, t, ', g, u, T, h, l, n, '}'

$\begin{array}{|c|c|c|c|c|} \hline \text{б} & \text{o} & \text{с} & \text{р} & \text{т} \\ \hline \end{array} \quad \begin{array}{|c|c|c|c|c|} \hline \text{г} & \text{o} & \text{с} & \text{р} & \text{т} \\ \hline \end{array}, \quad \begin{array}{|c|c|c|c|c|c|c|} \hline \text{г} & \text{o} & \text{с} & \text{р} & \text{o} & \text{т} & \text{o} & \text{у} \\ \hline \end{array}$
 $\begin{array}{|c|c|c|c|c|} \hline \text{с} & \text{o} & \text{с} & \text{р} & \text{т} \\ \hline \end{array} \quad \begin{array}{|c|c|c|c|c|} \hline \text{н} & \text{o} & \text{л} & \text{н} & \text{т} \\ \hline \end{array}.$

1) $\{\text{с}, \text{o}, \text{р}, \text{т}, ' ', \text{г}, \text{у}, \text{т}, \text{н}, \text{л}, \text{н}, '!', \text{ср}\}$

$\begin{array}{|c|c|c|c|c|} \hline \text{б} & \text{o} & \text{с} & \text{р} & \text{т} \\ \hline \end{array} \quad \begin{array}{|c|c|c|c|c|} \hline \text{г} & \text{o} & \text{с} & \text{р} & \text{т} \\ \hline \end{array}, \quad \begin{array}{|c|c|c|c|c|c|c|} \hline \text{г} & \text{o} & \text{с} & \text{р} & \text{o} & \text{т} & \text{o} & \text{у} \\ \hline \end{array}$
 $\begin{array}{|c|c|c|c|c|} \hline \text{с} & \text{o} & \text{с} & \text{р} & \text{т} \\ \hline \end{array} \quad \begin{array}{|c|c|c|c|c|} \hline \text{н} & \text{o} & \text{л} & \text{н} & \text{т} \\ \hline \end{array}.$

2) $\{\text{с}, \text{o}, \text{р}, \text{т}, ' ', \text{г}, \text{у}, \text{т}, \text{н}, \text{л}, \text{н}, '!', \text{ср}\}$
 обр

Embeddings

Word2Vec

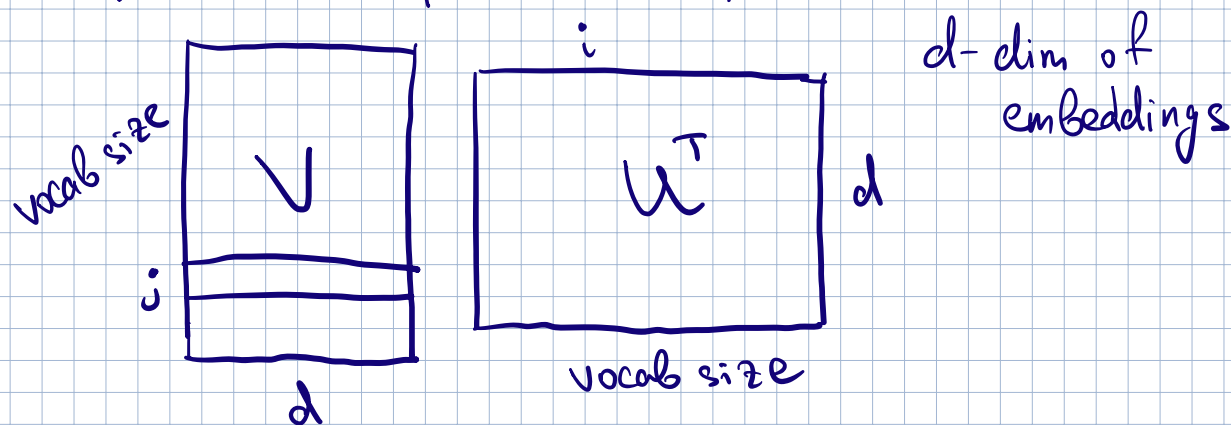
w_1 $\left[\begin{array}{c} w_2 \quad w_3 \quad \boxed{w_4} \quad w_5 \quad w_6 \end{array} \right] w_7 \dots$

- Continuous Bag of Words (CBOW)

$$p(w_4 | w_2, w_3, w_5, w_6)$$

- Skip-gram

$$p(w_2 | w_4), p(w_3 | w_4), p(w_5 | w_4), p(w_6 | w_4)$$



$$p(w_i | w_j) \approx \frac{e^{v_j^T u_i}}{\sum_{k=1}^{\text{vocab size}} e^{v_j^T u_k}} \quad \begin{matrix} v_j \in \mathbb{R}^d \\ u_i \in \mathbb{R}^d \end{matrix}$$

$$- \sum_{\substack{w_j - \text{center} \\ w_i - \text{context}}} \log p(w_i | w_j) \rightarrow \min_{u, v}$$

For CBOW $p(w_j | w_{j-n}, w_{j-n+1}, \dots, w_{j-1}, w_{j+1}, \dots, w_{j+n})$
 $(v_{j-n} + v_{j-n+1} + \dots + v_{j+n})^T u_j$

Multiclass \rightarrow Binary classification

$\left. \begin{matrix} w_i - \text{context word} \\ w_j - \text{center word} \end{matrix} \right\}$ "Do these words appear together"

$$p(w_i \text{ from context of } w_j) = \sigma(v_j^T u_i)$$

$$- \sum_{\substack{w_j - \text{center} \\ w_i - \text{context}}} \log \sigma(v_j^T u_i) \rightarrow \min_{u, v}$$

Negative sampling

$$- \sum_{\substack{w_j - \text{center} \\ w_i - \text{context}}} \log \sigma(v_j^T u_i) - \sum_{(w_j, w_i) - \text{neg. pair}} \log (1 - \sigma(v_j^T u_i)) \rightarrow \min_{u, v}$$

nn. Embedding index \rightarrow token

- V
 - U^T
 - $\frac{1}{2}(V + U^T)$
- $\text{man} \leftrightarrow \text{king}$
 $\text{woman} \leftrightarrow \text{queen}$
 $V_{\text{man}} - V_{\text{king}} \approx V_{\text{woman}} - V_{\text{queen}}$

FastText - extended Skip-gram

$\begin{bmatrix} \text{b o s p u} & \text{g o s p u}, & \text{g o s p o t o u} \\ \text{s o s p u} & \text{n o l n u} \end{bmatrix}$

$p(\text{s o s p u} \mid \text{g o s p u})$ - Skip-gram

$\text{g o s p u} \rightarrow \langle \text{g o s p u} \rangle, \langle \text{go}, \text{gos}, \text{osp}, \text{spu}, \text{pu} \rangle$

$p(\langle \text{s o s p u} \rangle \mid \langle \text{g o s p u} \rangle)$ $n=3$

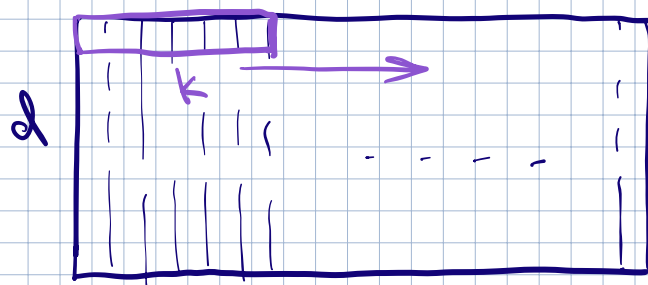
$p(\langle \text{s o s p u} \rangle \mid \langle \text{go} \rangle)$

$p(\langle \text{s o s p u} \rangle \mid \langle \text{gos} \rangle)$

...

GloVe - Global Vectors

Text CNN 1-D convolutions



d - channels
 L - similar to height and width

Conv. kernel $k \times 1 \times d \times d'$

$$d \times L \rightarrow d' \times L'$$

- Variable length

$$B \times C \times H \times W$$

