

① One hot encoding

[ One hot encoding is a method used to represent text (words / tokens) as binary vectors in NLP ]

Each unique word in the dataset is represented as a vector, where all elements are 0 except for single element that corresponds to that word, which is set to 1

Example dataset:

1. "cat eats fish"
2. "dog eats fish"
3. "dog chased cat"

① create a vocabulary:

["cat", "eats", "fish", "dog", "chases"]  
 unique words (vocabulary)

② assign index to each word

cat → 0

eats → 1

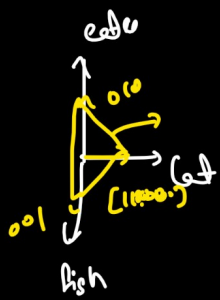
fish → 2

dog → 3

chases → 4

③ create one hot vectors

cat eats fish dog ch



"cat eats fish"

$$\begin{bmatrix} 1 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 0 & 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} 0 & 0 & 1 & 0 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 0 & 0 & 0 & 0 \end{bmatrix}$$

$$\begin{bmatrix} 0 & 1 & 0 & 0 & 0 \end{bmatrix}$$

$$\begin{bmatrix} 0 & 0 & 1 & 0 & 0 \end{bmatrix}$$

$$\begin{bmatrix} 0 & 0 & 0 & 1 & 0 \end{bmatrix}$$

$$\begin{bmatrix} 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

"dog eats fish"

$$\begin{bmatrix} 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 0 & 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} 0 & 0 & 1 & 0 \end{bmatrix}$$

"dog chews cat"

$$\begin{bmatrix} 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix} \rightarrow \text{sparse}$$

Advantages

① easy to implement with python

Disadvantages

① sparse matrix  $\rightarrow$  overfitting

② issues with size  
ML algorithms  $\rightarrow$  fixed size input

③ out of vocabulary

④ no semantic meaning is getting captured