

one hot encoding:

D-1 \rightarrow Tom purchased food.

D-2 \rightarrow cat eats food.

D-3 \rightarrow Dog eats food.

D-1 \rightarrow $\begin{bmatrix} [1, 0, 0, 0, 0, 0], \\ [0, 1, 0, 0, 0, 0], \\ [0, 0, 1, 0, 0, 0] \end{bmatrix}$

vocabulary: $\odot \cdot \times \cdot$

\rightarrow [Tom, purchased, food, cat, eats, dog]

D-2 \rightarrow $\begin{bmatrix} [0, 0, 0, 1, 0, 0], \\ [0, 0, 0, 0, 1, 0], \\ [0, 0, 1, 0, 0, 0] \end{bmatrix}$

D-3 \rightarrow $\begin{bmatrix} [0, 0, 0, 0, 0, 1], \\ [0, 0, 0, 0, 1, 0], \\ [0, 0, 1, 0, 0, 0] \end{bmatrix}$

Test data \Rightarrow parrot eats food; {parrot is not available
in training vocabulary}

Advantage \Rightarrow Easy implementation (pd.get_dummies, onehotencoder)

Disadvantage

\Rightarrow sparse matrix \Rightarrow [10k words in vocabulary]

\Rightarrow It won't perform with dynamic length document.

\Rightarrow Run out of vocabulary (OOV)

\Rightarrow No semantic meaning to the words

\Rightarrow This is good batch

$\begin{bmatrix} [1, 0, 0], [0, 1, 0], \\ [0, 0, 1] \end{bmatrix}$

