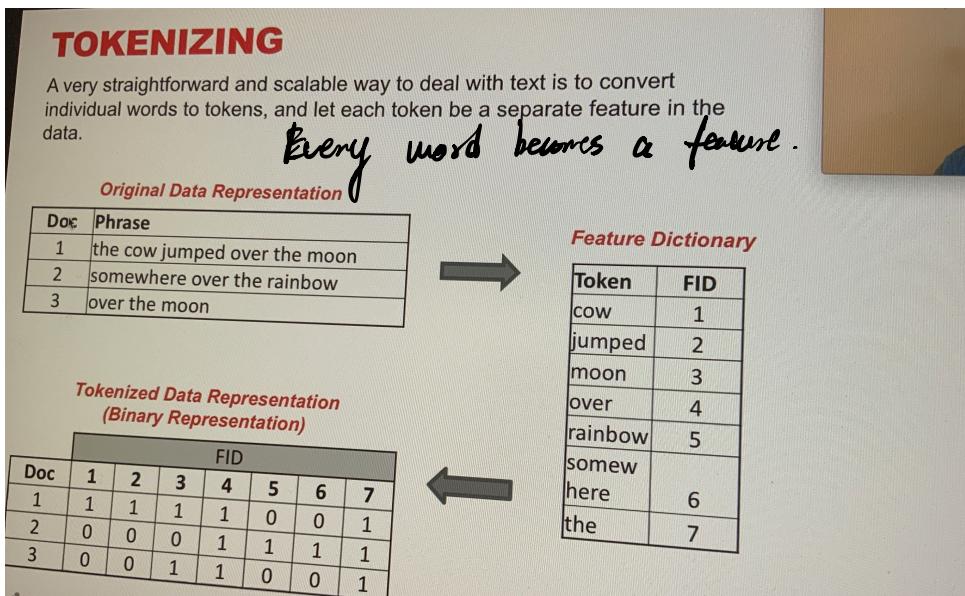


# Text Processing:

## Tokenization:



$$TF \cdot IDF = \text{Term Frequency} \times \text{Inverse Document Frequency}$$

TF measures how many times a particular token appears in a particular document, and gives an indication of how important that token is to the document.

$$TF(\text{term}, \text{Doc}) = \sum_{\text{word} \in \text{Doc}} I(\text{term} == \text{word})$$

IDF measures how often a word appears across all documents.

It penalizes those words that appear frequently

$$IDF(\text{term}, \text{Corpus}) = \log \frac{|\text{corpus}|}{\sum_{\text{doc} \in \text{corpus}} I(\text{term} \in \text{doc})}$$

## Word Embeddings:

Represent a word by a dense vector such that words with similar meaning or sentence are close to each other in the vector space.  
(measured by cosine similarity)

## Naive Bayes:

$$\text{posterior} = \frac{\text{prior} \times \text{likelihood}}{\text{evidence}}$$

$$\begin{aligned} P(y|x_1, x_2, \dots, x_m) &= \frac{P(y) P(x_1, x_2, \dots, x_m | y)}{P(x_1, x_2, \dots, x_m)} \\ &= \frac{P(y) \cdot P(x_1|y) P(x_2|y) \dots P(x_m|y)}{P(x_1, x_2, \dots, x_m)} \\ &= \frac{P(y) \cdot \prod_{i=1}^m P(x_i|y)}{P(x_1, \dots, x_m)} \end{aligned}$$

$$\hat{y} = \underset{y \in Y}{\operatorname{argmax}} P(y) \cdot \prod_{i=1}^m P(x_i|y)$$

if binary, this will need  $2n$  parameters.  
+ 2 prior parameters.

## Spam Detection:

Deliberately put some spam emails to inbox to avoid the negative feedback loops caused by selection bias.

ex of selection bias  $\rightarrow$  negative feedback loop:

- If 60% of African princess related email go to the spam email, we will never learn the future. when we train our model because we are using user's feedback on Inbox email.