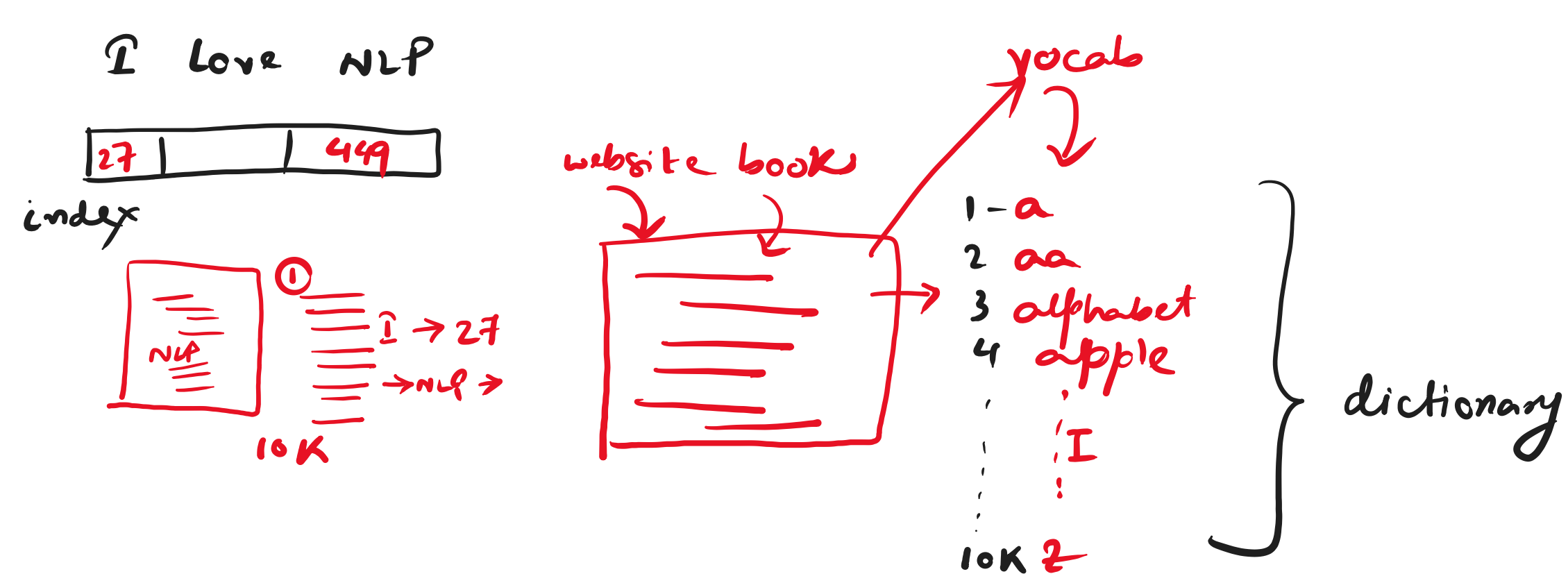
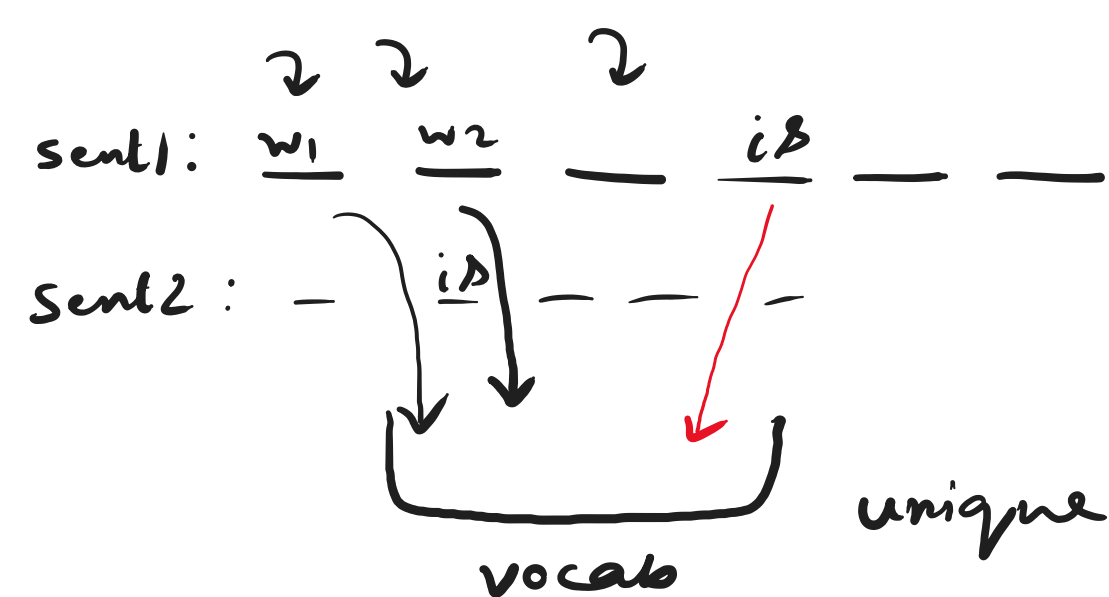


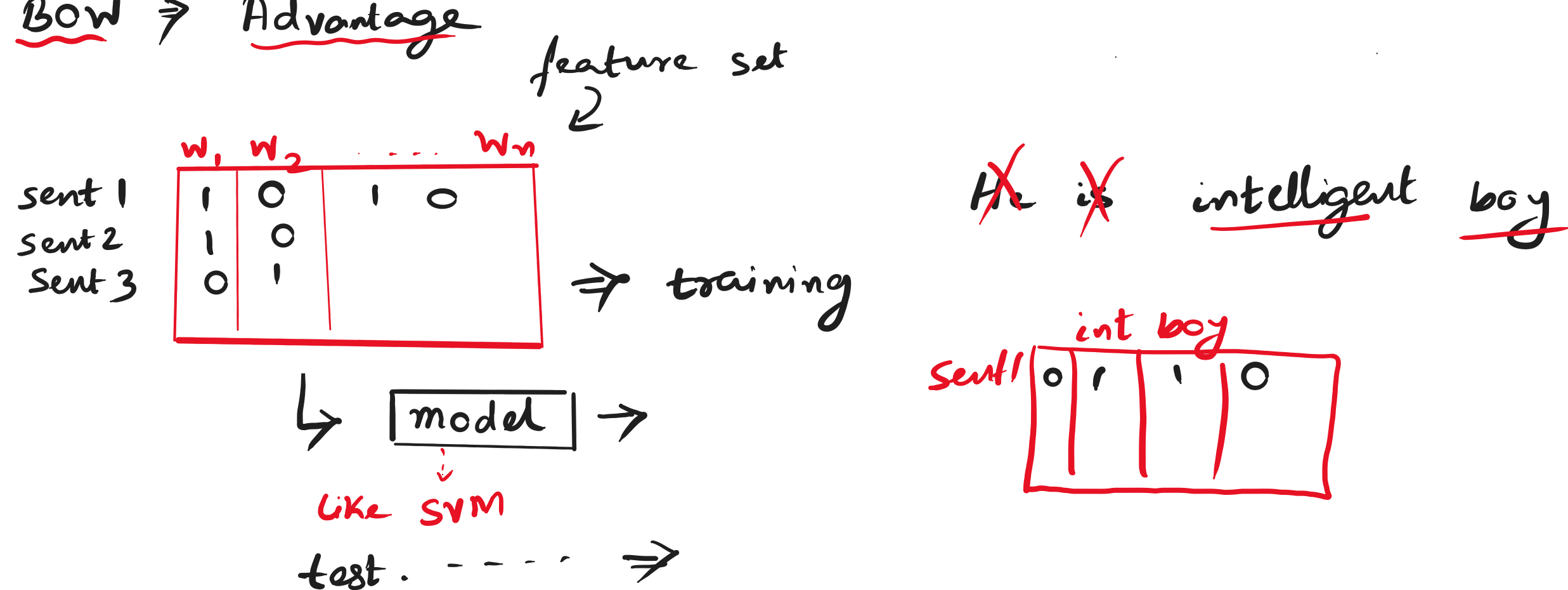
## Recap

- Text  $\rightarrow$  numerical representation of BOW
- Coding implementation

Topic  $\rightarrow$  BOW, TF-IDF, Coding



## BOW $\rightarrow$ Advantage



- Easiest
- Small dataset, BOW works well
- if " is large, BOW " X

## Disadvantages

- Big dataset, X
- sparsity
- semantic information X
- Context X
- Order X
- Word weightage is not give

	$w_1$	$(w_2)$	$(w_3)$	$w_4$
sent1	0	intelligent	boy	0
		0.1	0.9	

★ imp

## TF - IDF (Term Frequency - Inverse Document Frequency)

sent1: intelligent boy  
sent2: intelligent girl  
sent3: intelligent boy girl

$$\text{Term Frequency (TF)} = \frac{\text{No. of rep. of words in sentence}}{\text{Total no. of words in sentence}}$$

$$\text{IDF} = \log \left( \frac{\text{Total no. of sentences}}{\text{No. of sentences containing words}} \right)$$

$$\text{TF} * \text{IDF} = \checkmark$$

histogram	words		TF				IDF	
	words	Frequency	words	sent1	sent2	sent3	words	IDF
$\left\{ \begin{array}{l} \text{intelligent} \\ \text{boy} \\ \text{girl} \end{array} \right.$	intelligent	3	intelligent	1/2	1/2	1/3	intelligent	$\log(3/3) = 0$
	boy	2	boy	1/2	0	1/3	boy	$\log(3/2) =$
	girl	2	girl	0	1/2	1/3	girl	$\log(3/2) =$

	$f_1$	$f_2$	$f_3$	
	intelligent	boy	girl	
sent1	0	$1/2 * \log(3/2)$	0	$\rightarrow$ unique words (feature)
sent2	0	0	$1/2 * \log(3/2)$	
sent3	0			
				$\rightarrow$ feature set
				$\Rightarrow$ training data