

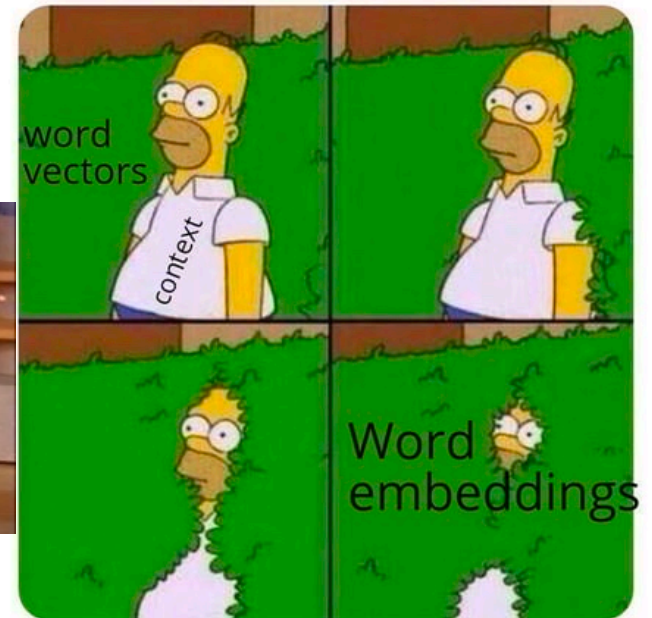
May 1, 2023.

DSML: NLP module.

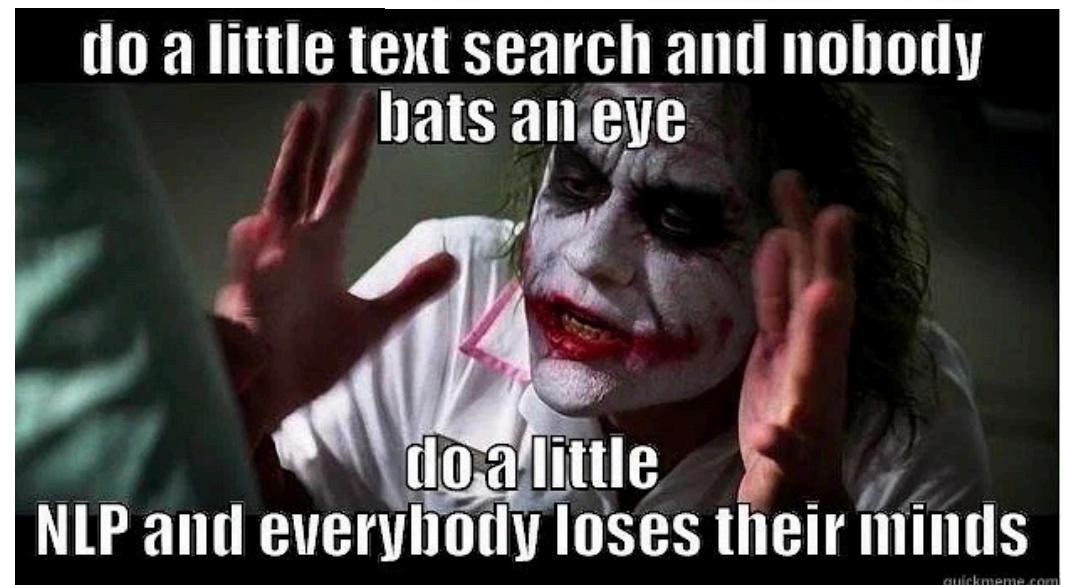
Word embeddings
in a nutshell

Topic Modeling.

Class starts
@ 9:05



When you penalize your Natural Language Generation model for large sentence lengths



Recap:

* Document Vectorization: Bag of Words (BoW)
TF/IDF.

* Word Vectorization: Continuous BoW,
Skip-gram.

* Language Modeling: Naïve Bayes to predict
the next word.

Agenda:

Classical NLP (No deep learning).

→ Parts-of-Speech tagging.

* Naïve approach, using heuristics and our knowledge of English grammar.

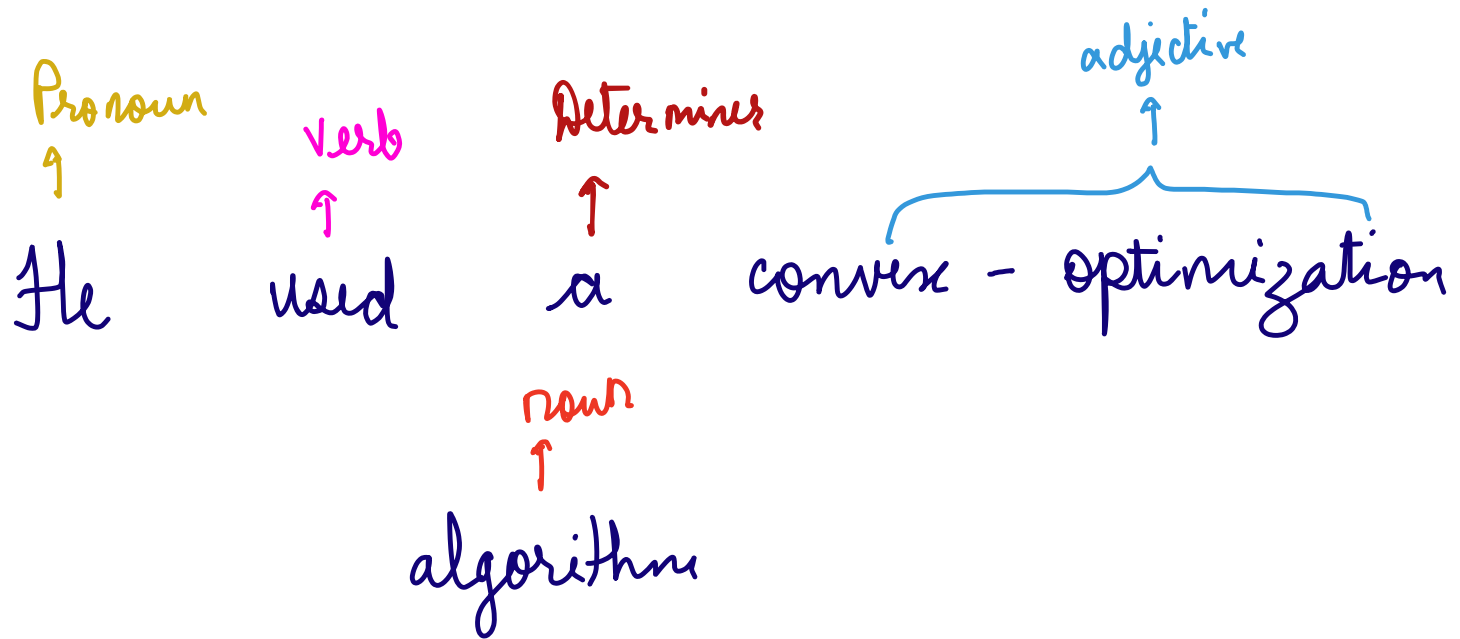
→ Topic Modeling.

* Statistical approach to build a generative model.

* The approach is called Latent Dirichlet Allocation (LDA).

* Business Case: Extract "topics" of discussion.

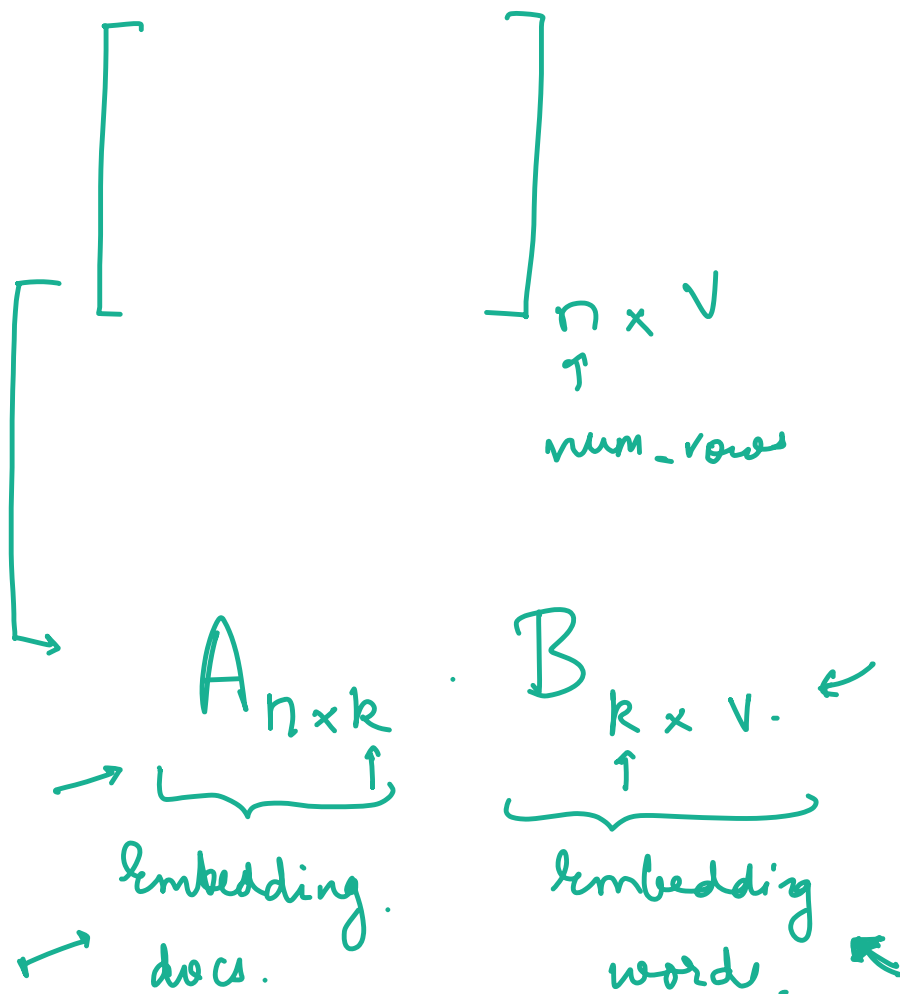
Parts of speech tagging



Hidden Markov models.

Topic Modeling.

→ Matrix factorization.



Corpus.



unique words

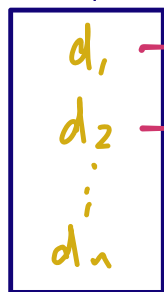


Vocabulary.

$V \rightarrow$ size of Vocabulary

Topic modeling as Clustering.

Corpus.

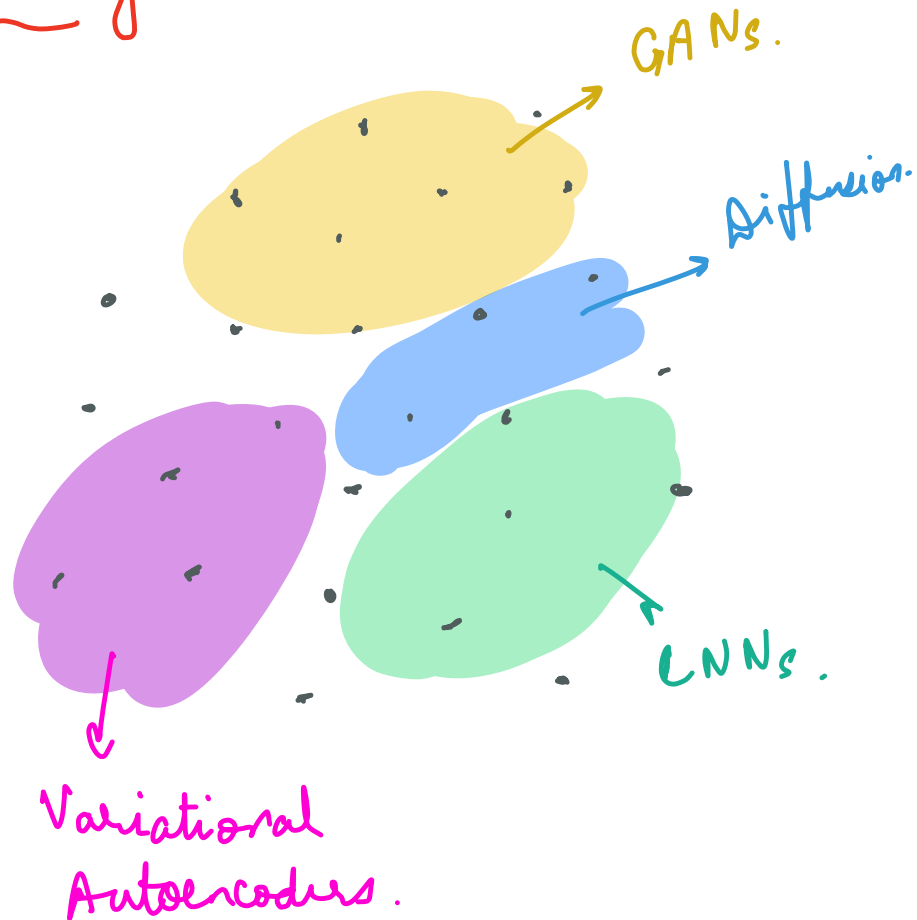


unique words



Vocabulary.

$V \rightarrow$ size of Vocabulary



Latent Dirichlet Allocation:

- k : No. of topics.
- V : Vocab size.
- M : No. of documents.
- N : No. of words in each document

z_{ij} K $[:::]$ V θ_{ij}

(2) A.P.D

2. Based on the ball you pick, you're sent to another ground "Beta"

1. Pick a ball from the ground "Theta"

M K $[0 0 0 0]$
 $[0 0 0 0]$

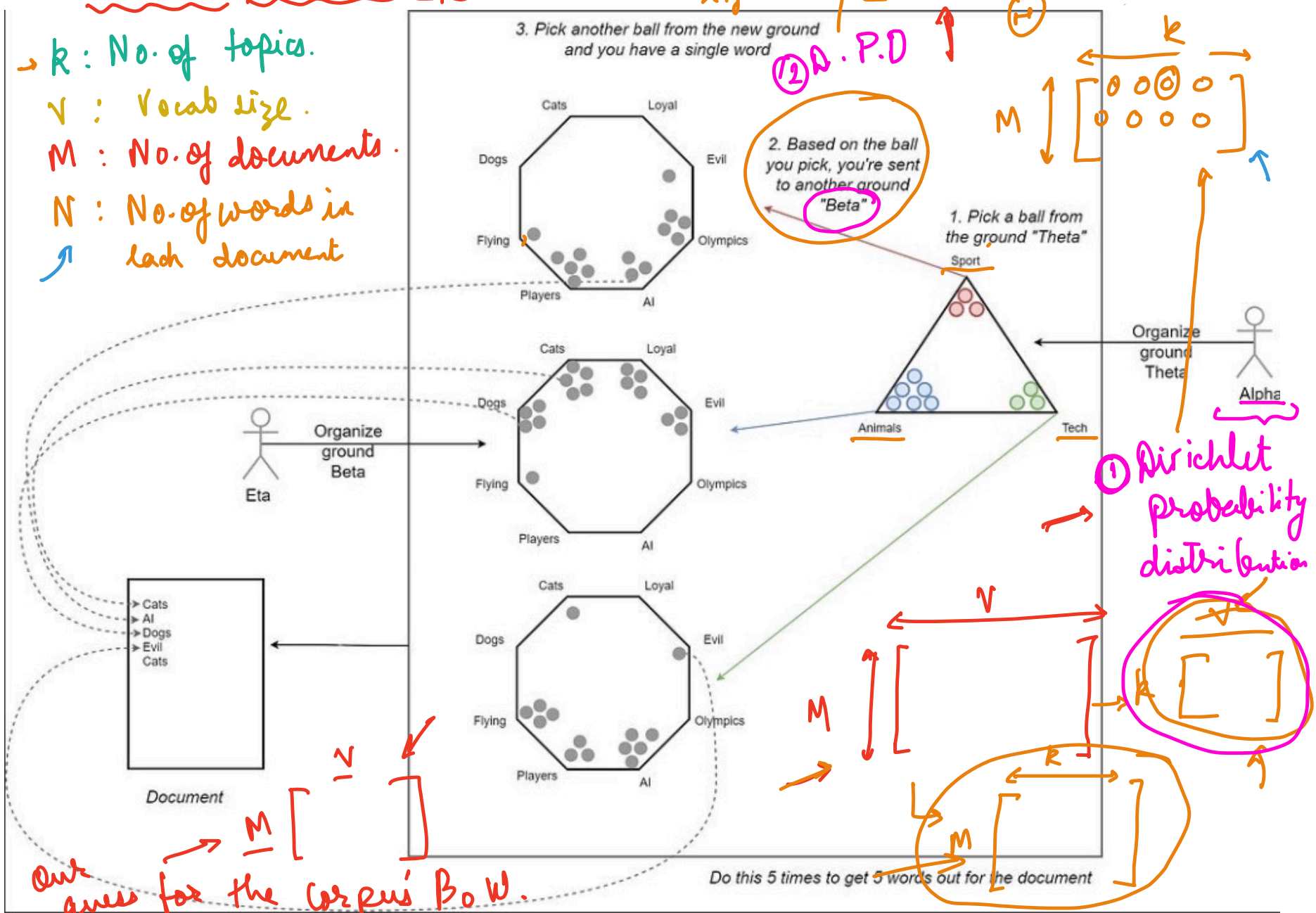
(1) Dirichlet probability distribution

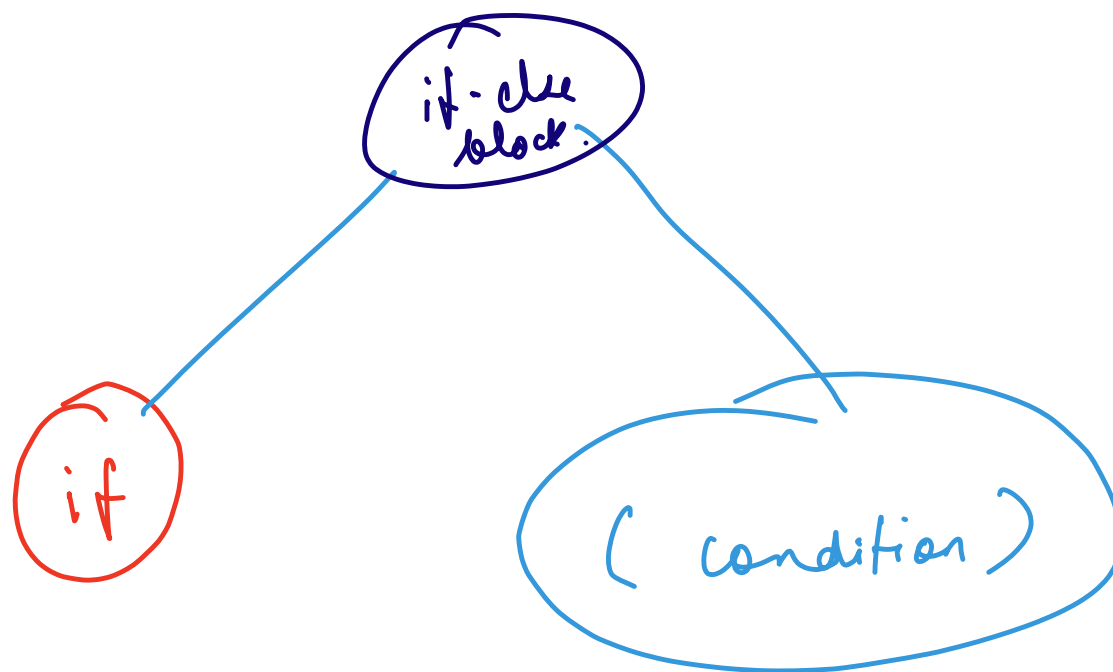
K $[]$

M K $[]$

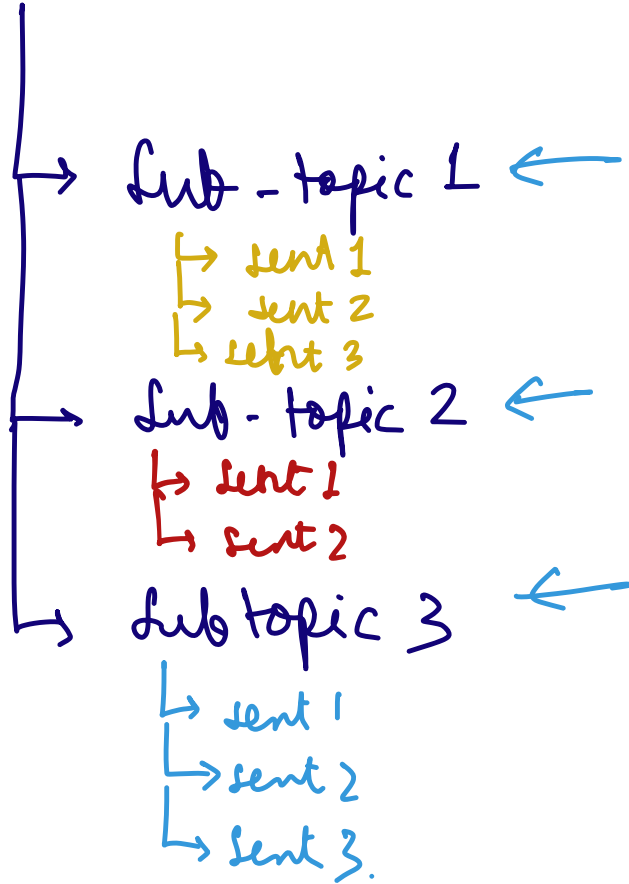
Do this 5 times to get 5 words out for the document

our guess for the corpus BOW.





(Main) Topic

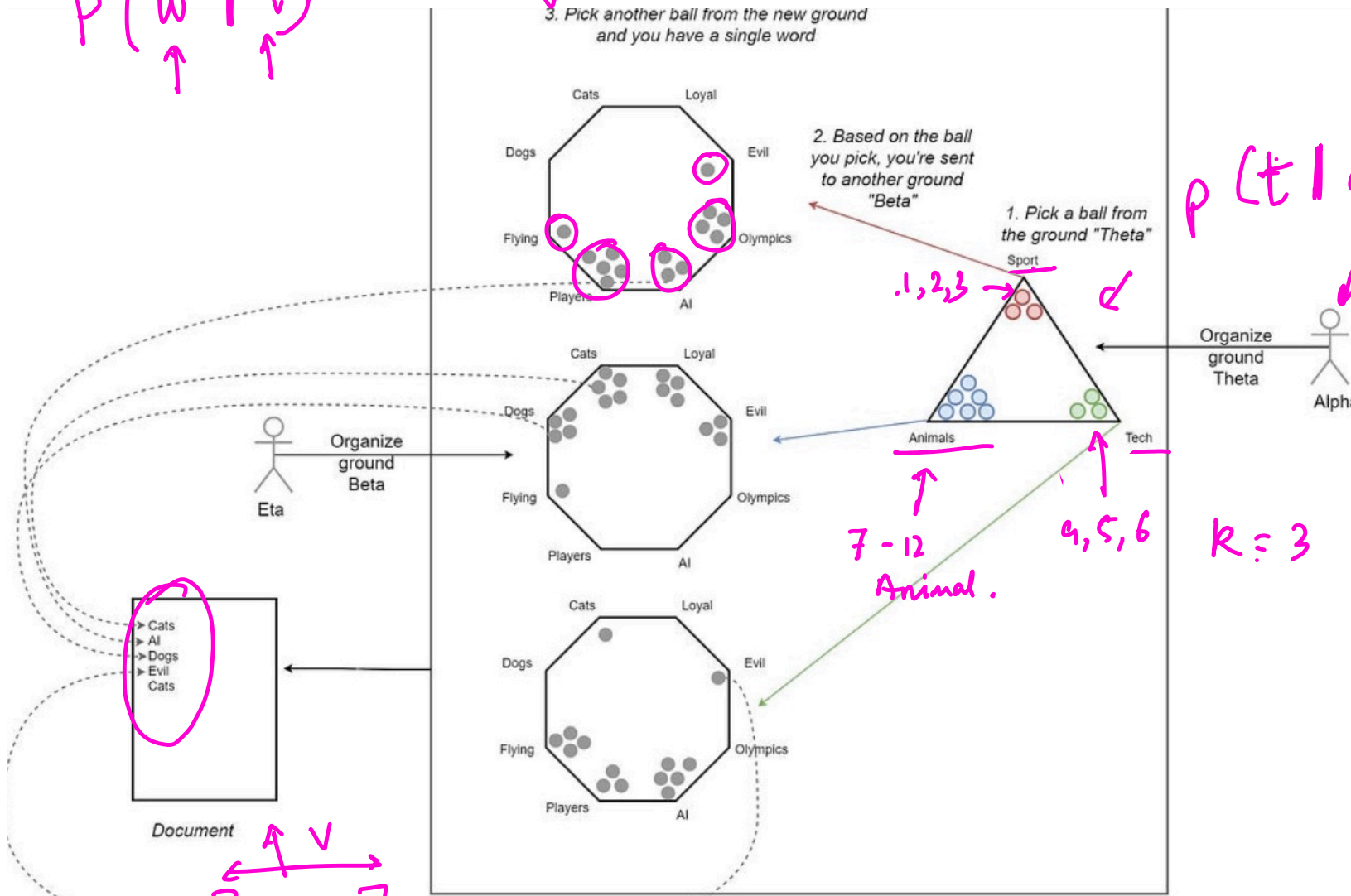


$$P(w | t)$$

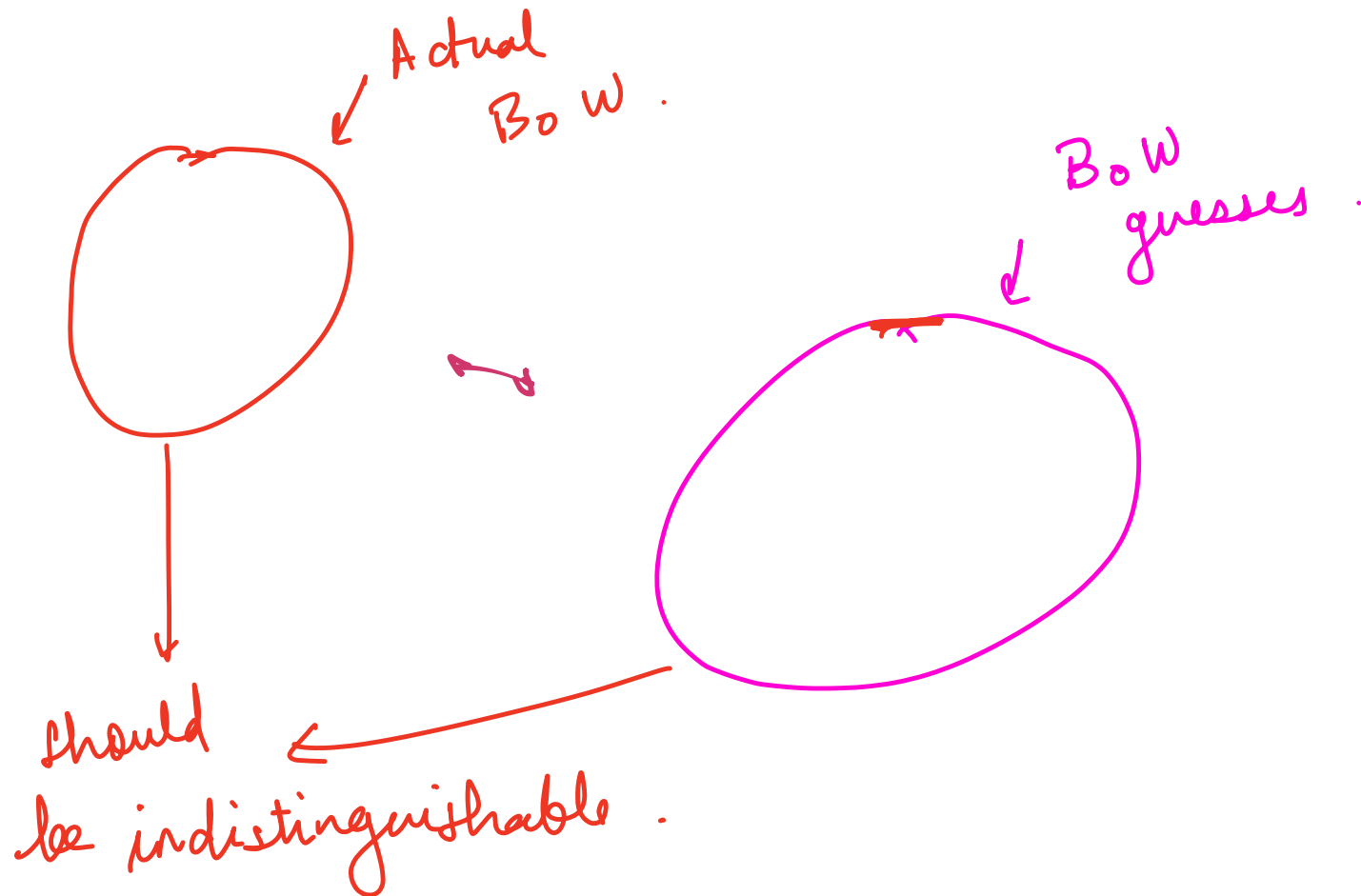
$$k \left[\frac{1}{14} \quad \frac{4}{14} \quad \frac{3}{14} \dots \right]$$

$$P(t | d)$$

$$k = 3$$



Now how to do topic modeling?



loss function: KL - Divergence.

```
[(0,  
  '0.020*"great" + 0.014*"good" + 0.012*"sound" + 0.012*"strings" + 0.011*"quality" + 0.011*"price"  
+ 0.010*"would" + 0.009*"time" + 0.008*"use" + 0.007*"well"'),
```

Not much to write about here, but it does exactly what it's supposed to. filters out the pop sounds. now my recordings are much more crisp. it is one of the lowest prices pop filters on amazon so might as well buy it, they honestly work the same despite their pricing,



Matrix 1:

$$\Theta = \begin{bmatrix} \theta_{ij} \end{bmatrix}$$

$M (\text{Num_docs})$

$k (\text{num_topics})$

$\theta_{ij} \rightarrow$ Probability that topic j is important for document M

$$P(t | d)$$

Matrix 2:

$$Z = \begin{bmatrix} z_{ij} \end{bmatrix}$$

$\xleftarrow{\quad} \underset{V}{\text{(vocab size)}} \xrightarrow{\quad}$

$\xleftarrow{\quad} k \text{ (num-topics)} \xrightarrow{\quad}$

$z_{ij} \rightarrow$ Probability that word j is relevant to topic i .

$$P(w | t)$$

$$\rightarrow P(w|d) = \frac{P(w \cap d)}{P(d)}$$

$$P(w|t)$$

$$\frac{P(w \cap d \cap t)}{P(w \cap d \cap t)}$$

$$= P(w \cap d | t) \cdot P(t)$$

$$P(t|d)$$

$$P(w|t) = \frac{P(w \cap t)}{P(t)}$$

$$P(t|d) = \frac{P(t \cap d)}{P(d)}$$

$$P(w|t) \cdot P(d)$$

$$\rightarrow P(w|t)$$

J.P:

$$P(w, d | t) P(t)$$

$$= P(w|t) \cdot P(t)$$

$$\uparrow P(d)$$