**\* Topic Modelling**

$K \longrightarrow$ # of Topics
$N \longrightarrow$ # of words
$C_m \longrightarrow$ Distribution over topics
$K \longrightarrow$ index $[1, \ldots, K]$
$w \longrightarrow$ index words $[1 \ldots |V|]$

} PAPER :
David Blei, Jordan
& Ng 2003

**\* Word embedding**

Dictionary $\longrightarrow V$
# of words $\longrightarrow |V|$

] $\longrightarrow$ Glove        Stanford } two famous
$\longrightarrow$ word2vec  Google } word embeddings

word2vec $\Big\langle$ 
CBow (continuous bag of word) (Dealing row (w) in co-occurrence table)
skipgram (dealing column in co-occurrence table)

**\* Tagging** $\longrightarrow$ 
PoS
NER (People, locations, organizations & misc)
CHONKING (entity detection)
NER (entity classification)

**\* SEQ to SEQ**

| Text | Image |
|---|---|
| <u>Text</u> | <u>Image</u> |

<u>Text</u>

* Corpus / Collection
* Document
* Word (like pixels / superpixels)
* Vocabulary

<u>Image</u>

* Training Set $\longrightarrow S$
* Total images in $S \longrightarrow N$
* For each image, $x$:
  $\longrightarrow$ Let "$T$" be its total superpixels

So, $x = [x^1, x^2, \cdots, x^T]$

$\underbrace{\qquad\qquad}_{T \text{ superpixels}}$

And now:

$x^1 \in \mathbb{R}^n$  ("$n$" features for each superpixel)

$x \in \mathbb{R}^D$

NLP Applications

- → Document Classification (Supervised)
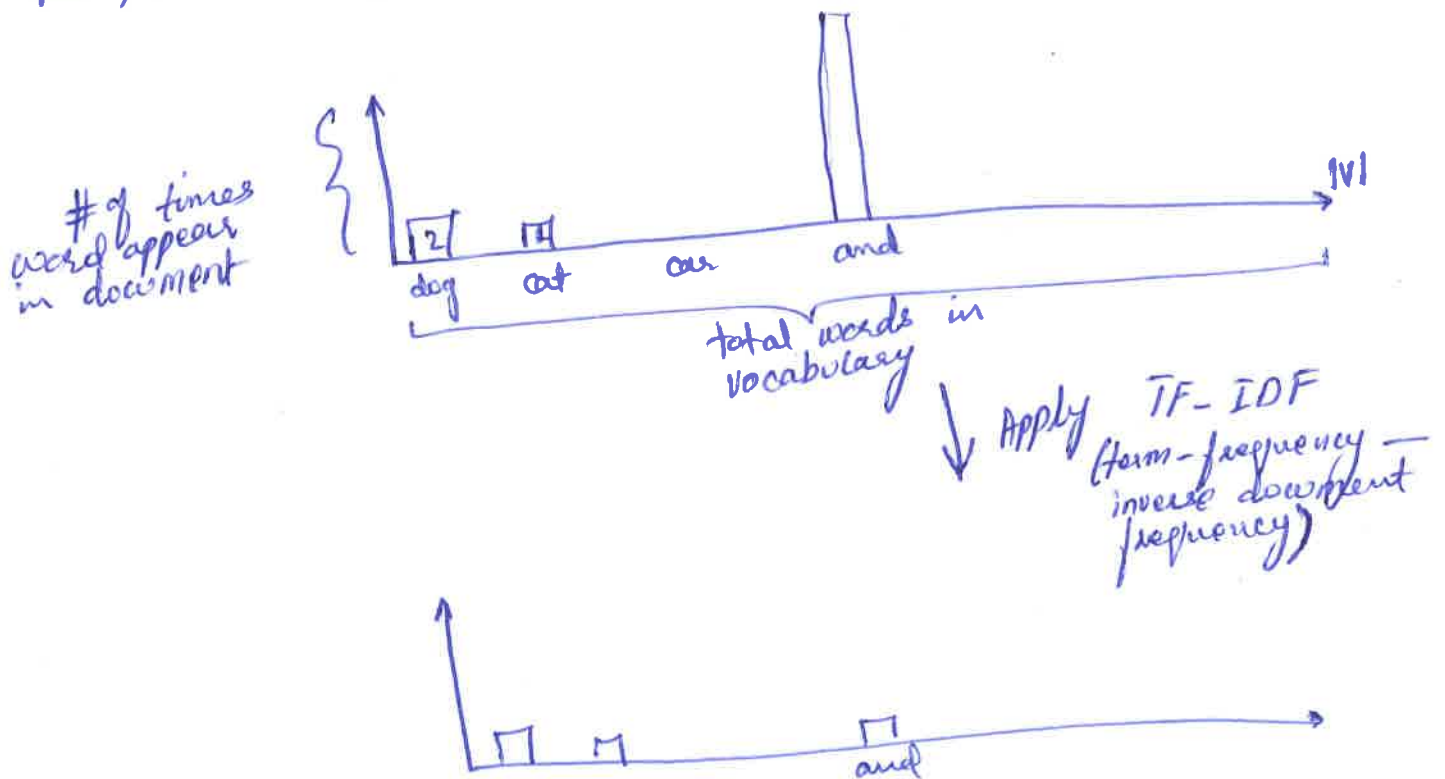- → Topic Modelling (Unsupervised)
- → word embedding
- → Tagging
- → Seq2Seq

# Document Classification

Assign a document its label

$$X \in R^D \longrightarrow y$$

where $X$ is document

Now, how to represent document as a vector? (As histogram of words)



# of times word appear in document

dog   cat   car   and

total words in vocabulary

Apply TF-IDF (term-frequency - inverse document frequency)



and

TF-IDF: statistical measure that tells how important a word is in the document.

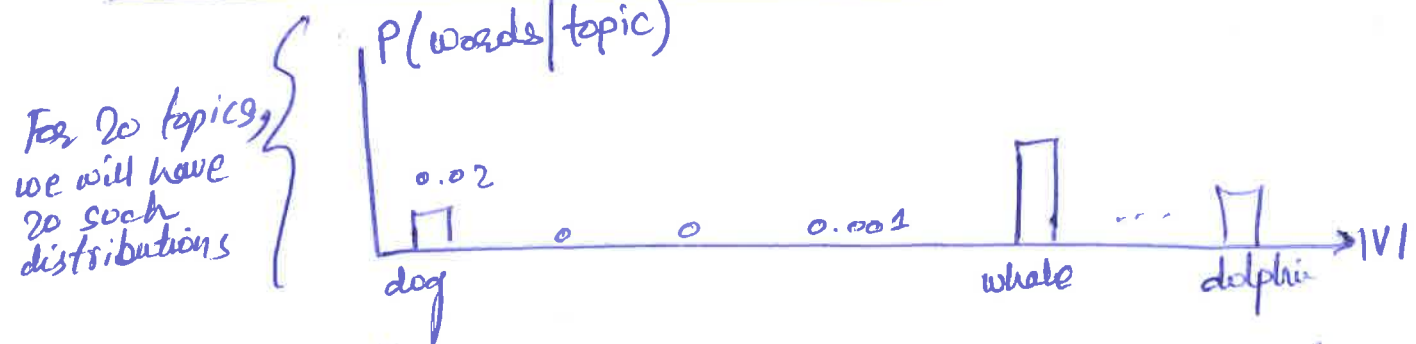# Topic Modelling (Unsupervised Problem)

### Example

Find 20 main topics in

a corpus of "N" documents

collection OR Training Set

like "N" images

### Topic

A topic is a distribution over the words in Vocabulary "V"

$$|V| \sim 50,000 - 200,000$$

For 20 topics, we will have 20 such distributions

P(words|topic)

0.02

0

0

0.001

dog

whale

dolphin

→|V|

This vector can be taken for a problem of topic labelling

### PAPer/Code:

1) To understand topic modelling, read Latent Dirichlet Allocation (David Blei 2003)

2) Use scikit learn (for lda)

# Word Embeddings

- → Two famous word embeddings
  - → Glove (Stanford)
  - → Word2vec (Google)

## Example

Today, I parked my [car] on street.

<u>window = 3 words</u>

Column

words



|       | parked ... | my | on - - - - - |V| |
|-------|-----------|-----|-----|
| flag  |           |     |     |
| ⋮     |           |     |     |
| car   |           | +1  | +1  |
| ⋮     |           |     |     |
| truck |           |     |     |

|V|

Co-occurence matrix

We can find $P(w, C)$, $P(w|C)$ & $P(C|w)$ in co-occurence matrix

* In    word2vec
  - → CBOW (continuous bag of words)
    "Here we deal with rows in co-occurence matrix"
  - → Skipgram
    ("Here we deal with column in co-occurence matrix")

# TAGGING

Tag a word in a sentence

## POS tagging : { Noun, adj, verb... }

→ Old way of tagging : look up tables

→ New way :

    1) word embedding ($\sim 300 - 500D$)

    3) trigram ⟶ $900D$ [ concatenation of 3-words ] of word embeddings

    4) train classifier

$$900D \longrightarrow y$$
$$X$$

## NAMED - Entity Recognition (NER)

{ People
  location
  organisation
    MUC ⟶ events
      ⟶ number
      ⟶ time expression

Example    B-loc   I-loc
            LAS Vegas

                                I-org  I-org  I-org
            B-loc
            Commonwealth bank of Australia
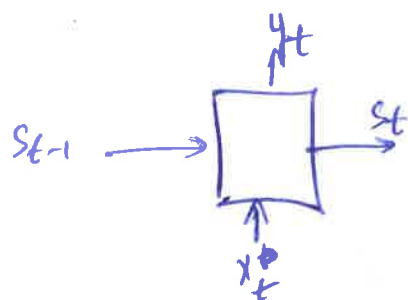
## CHUNKING    Entity (detection) or segmentation

    .... _entity_   ....   _entity_  .....
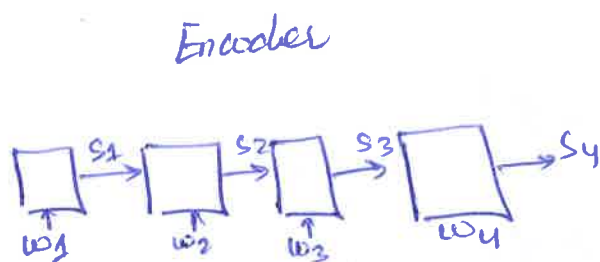
P: input
Q: output

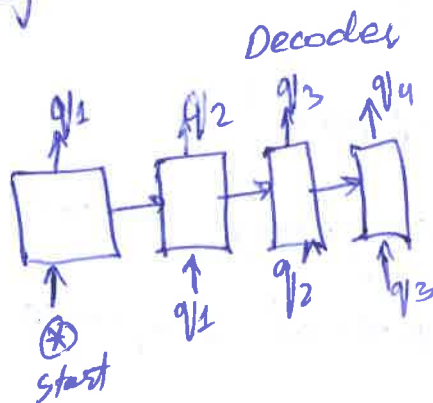

Three famous RNNs {
Jordan
Elman
LSTM
}

Generation of Q from P using Encoder-Decoder:

Encoder                                    Decoder



In attention Networks

$$S_1 \longrightarrow S_1 \times a_1$$
$$S_2 \longrightarrow S_2 \times a_2$$
$$S_3 \longrightarrow S_3 \times a_3$$

$a_1, a_2, a_3$ are attention weights. They can be learnt
as follows:
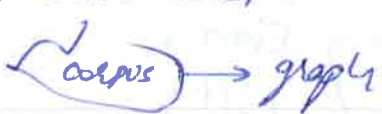


PAPer

1) Seq2Seq → SutsRever 2015
2) Seq2Seq Attention → Badhanan

# Other tasks

1) Co-reference Resolution:
   Shaukat did same ... He told

2) Relation Extraction:
   Father-son, mother-son ...

3) Taxonomy extraction,
   corpus → graph

4) Machine translation
   French → English → English improved

5) Summarization

6) Matching (BiMPM)

7) RTE (recognising textual entailment)

8) NLI (Natural language inference)