Word2Vec $\longrightarrow$ Tomas Mikolov $\rightarrow$ 2013 $\rightarrow$ Google

$\Downarrow$

Community support

Fast text $\longleftarrow$ Facebook
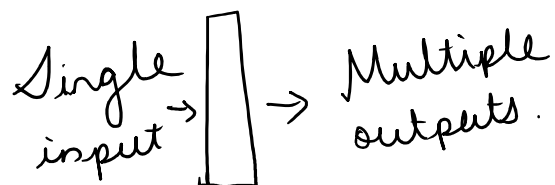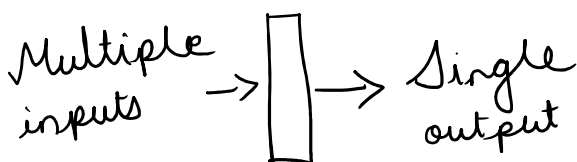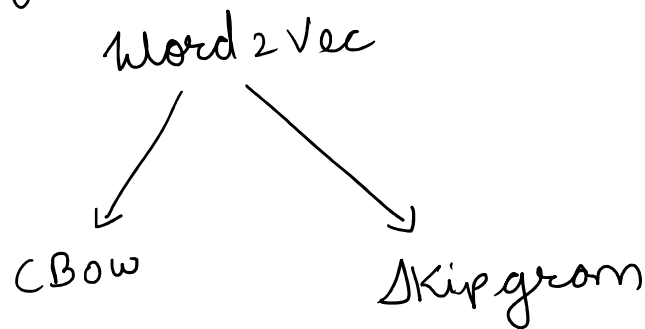
$\Downarrow$

Supports multiple languages

We con use Word2Vec in 2 ways:

① Pretrained model $\longrightarrow$ Download $\rightarrow$ use

$\rightarrow$ 3 Billions words $\rightarrow$ Google news dataset

② Customised model $\rightarrow$

Architectures of Word2Vec :$\rightarrow$

Word2Vec

CBow          Skipgram

Multiple inputs $\rightarrow$ ▯ $\rightarrow$ Single output          Single input $\rightarrow$ ▯ $\rightarrow$ Multiple outputs.
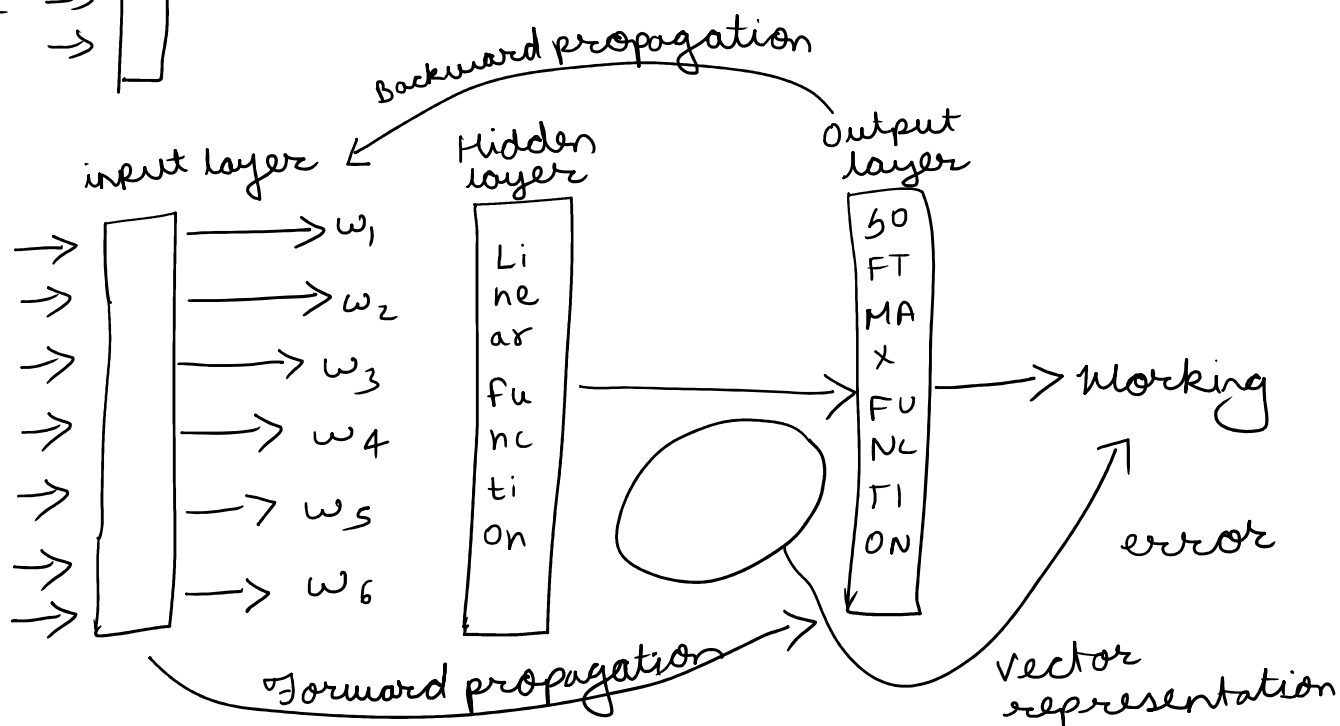
CBow $\rightarrow$ Continuous Bag of words

"Rajesh is [working] in MNC, he is a hardworking guy!"
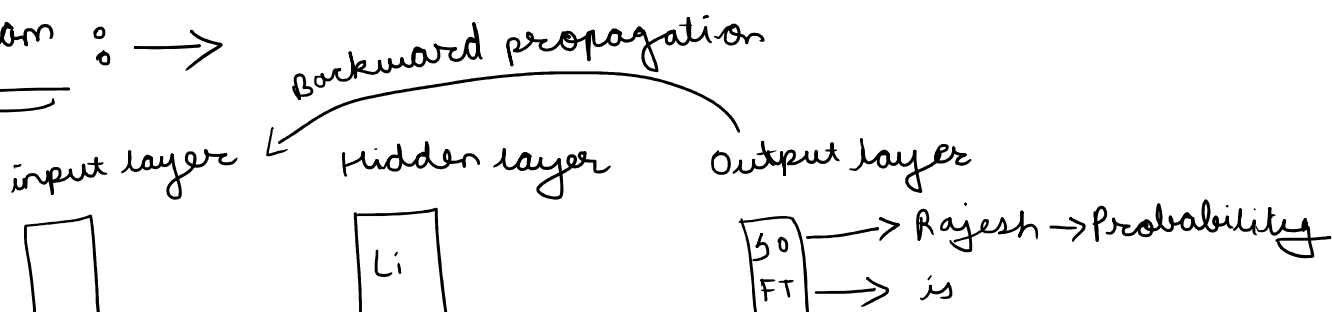
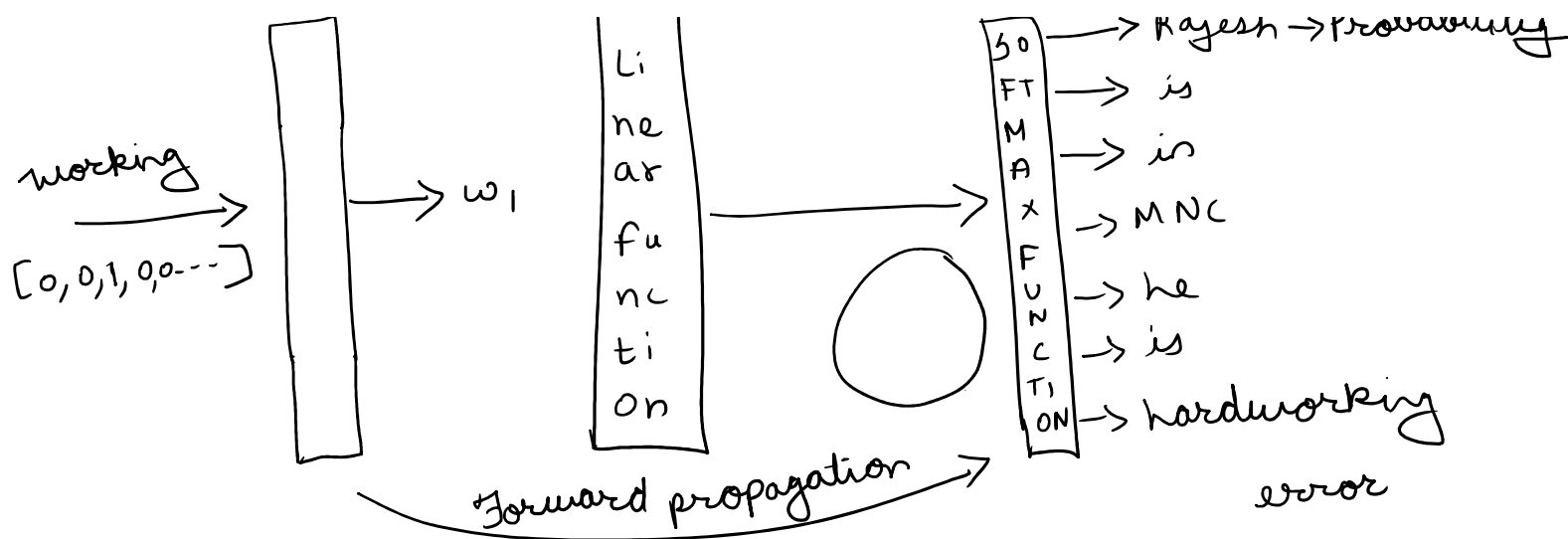Target word $\quad\quad\quad\quad$ contextual words

input layer

Rajesh $\rightarrow$
is $\rightarrow$
in $\rightarrow$
MNC $\rightarrow$
he $\rightarrow$
is $\rightarrow$

Rajesh $= [1, 0, 0, 0, - - - -]$

is $= [0, 1, 0, 0, - - -]$

in $= [0, 0, 0, 1, 0, - - - -]$

Backward propagation

input layer $\quad\quad$ Hidden layer $\quad\quad$ Output layer

$w_1$
$w_2$
$w_3$
$w_4$
$w_5$
$w_6$

Linear function $\quad\quad$ SOFTMAX FUNCTION $\rightarrow$ working

Forward propagation

error

vector representation

① Kipgram : $\rightarrow$

Backward propagation

input layer $\quad\quad$ Hidden layer $\quad\quad$ Output layer

Li

SOFT $\rightarrow$ Rajesh $\rightarrow$ Probability

$\rightarrow$ is

working

[0, 0, 1, 0, 0 ...] → □ |→ $w_1$ → Linear function → □ ○ → SOFTMAX FUNCTION → Rajesh → probability
→ is
→ in
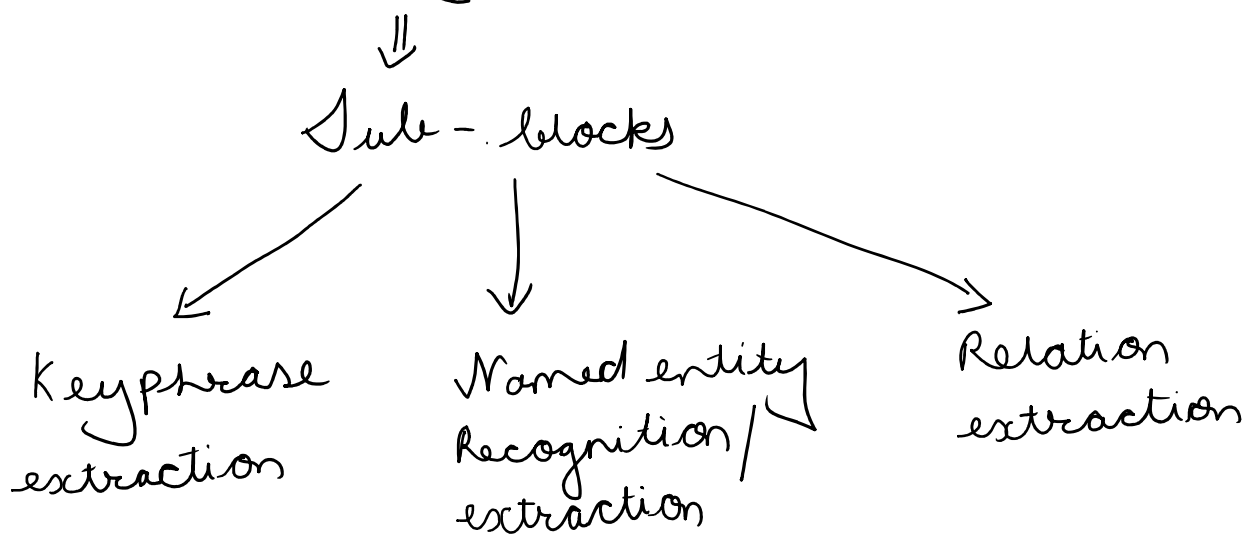→ MNC
→ he
→ is
→ hardworking

error

Forward propagation

Vector Representation of word = Weight Matrix × Probability of word.

## Information extraction : Big umbrella

→ we are extracting knowledge from the text.

$\|$

Sub - blocks

Keyphrase extraction

Named entity Recognition / extraction

Relation extraction

## Relation extraction: Open NRE

[Rajesh] was born in 1940, the only child of teacher [Aakash] and his wife [Rashi], a doctor--

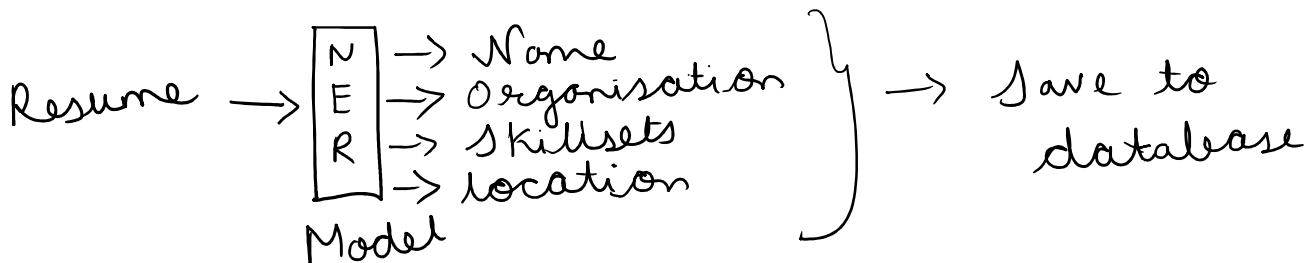teacher Aakash and his wife Rashi, a doctor..

Named entity Recognition / extraction ⇒ Spacy

⇓

Resume Parsing

Resume → [N E R] → Name
→ Organisation
→ Skillsets
→ location
Model
⎫
⎬
⎭ → Save to database

Keyphrase extraction:

→ extracting important phrases from the data.

ways to do Keyphrase extraction
↳ ① Ngrams
↳ ② Unsupervised algorithms
- RAKE, YAKE

Flow chart of unsupervised algorithm for Keyphrase extraction:

Keyphrase extraction.

Raw text $\longrightarrow$ Candidate generation

$\downarrow$

Candidate scoring

Final Ranking $\longleftarrow$ Post Processing $\longleftarrow$

① Candidate generation
Extract all phrases

② Candidate scoring
give scores to phrases
Important phrases

③ Post processing $\Rightarrow$ "game" $\longrightarrow$ "games"
insert "s"

SUN $\longrightarrow$ RAN

S $\rightarrow$ delete
R $\rightarrow$ insert $\Longrightarrow$ RAN
U $\rightarrow$ delete
A $\rightarrow$ insert

Laveshtein distance $\rightarrow$ game $\rightarrow$ games $\rightarrow$ 1
L.D. $\rightarrow$ 1

SUN $\longrightarrow$ RAN $\longrightarrow$ 4

Use of Laveshtein distance :

Single Review $\longrightarrow$ 1000 — 5000 words

Key phrases $\longrightarrow$ 2500

10 Documents $\longrightarrow$ 1,00,000 — 5,00,000 words

Key phrases $\longrightarrow$ 2,50,000

Important Keyphrases $\longrightarrow$ Keyphrases count $\downarrow$

Nice Movie
Nicely Movie $\Big]\longrightarrow 2$
$\Downarrow$

Final Ranking: $\longrightarrow$ 5,00,000 Phrases
$\Downarrow$
1,50,000 Phrases
$\Downarrow$
Top 100,500

Applications of Key phrase extraction : $\longrightarrow$

① Text summarization ( extractive )
② Data analysis
③ Root cause analysis

④ Chatbot

Hi, How are you? $\Longrightarrow$ Hello, How are you?