# Convert text to Numeric Vectors.

1. Bag of Words
2. Binary BoW  } n-grams
3. TF-IDF ✓
4. Word2Vec

$n=1$   unigrams
$n=2$   bi-grams
$n<3$   tri-grams.

crickets runs score ✓ ____

→ D1 : I love pasta. Pasta is tasty ✓→ Bigram

D2 : I hate pasta. Pasta is cheap ✗

→ Trigrams

global warming, Glaciers.

Size of the Vectors ⤳ No. of features.

(Vocab): list of unique words    8

Vocab :    I    love    pasta    hate    is    tasty    cheap ⤎ 7
                                                                    ?

|       | I | love | pasta | hate | is | tasty | cheap |
|-------|---|------|-------|------|-----|-------|-------|
| D1 :  | 1 | 1    | 2     | 0    | 1   | 1     | 0     |
| BBoW  |   | 1    | 1     | 0    | 1   | 1     | 1     |
| D2 :  | 1 | 0    | 2     | 1    | 1   | 0     | 1     |

→ Vocab Size

| | $x_1$ | $x_2$ | $x_3$ | $\cdots$ | $x_7$ | $\cdots$ | $x_{7800}$ | $y$ | → |
|------|----|----|----|----|----|----|----|----|---|
| D1 | 1 | 1 | 2 | $\cdots$ | 0 | | | | |
| D2 | 1 | 0 | 2 | 1 | $\cdot$ $\uparrow$ | | | +ve | |

→ Classification / clustering

## TF - IDF

→ Term Frequency – Inverse Document Frequency

$$\text{Term Frequency (TF)} = \frac{\text{No. of times the word occurs in the doc.}}{\text{Total no. of words in the doc.}}$$

$$\text{Document Frequency (DF)} = \frac{\text{No. of documents in which the word occurs.}}{\text{Total no. of documents.}}$$
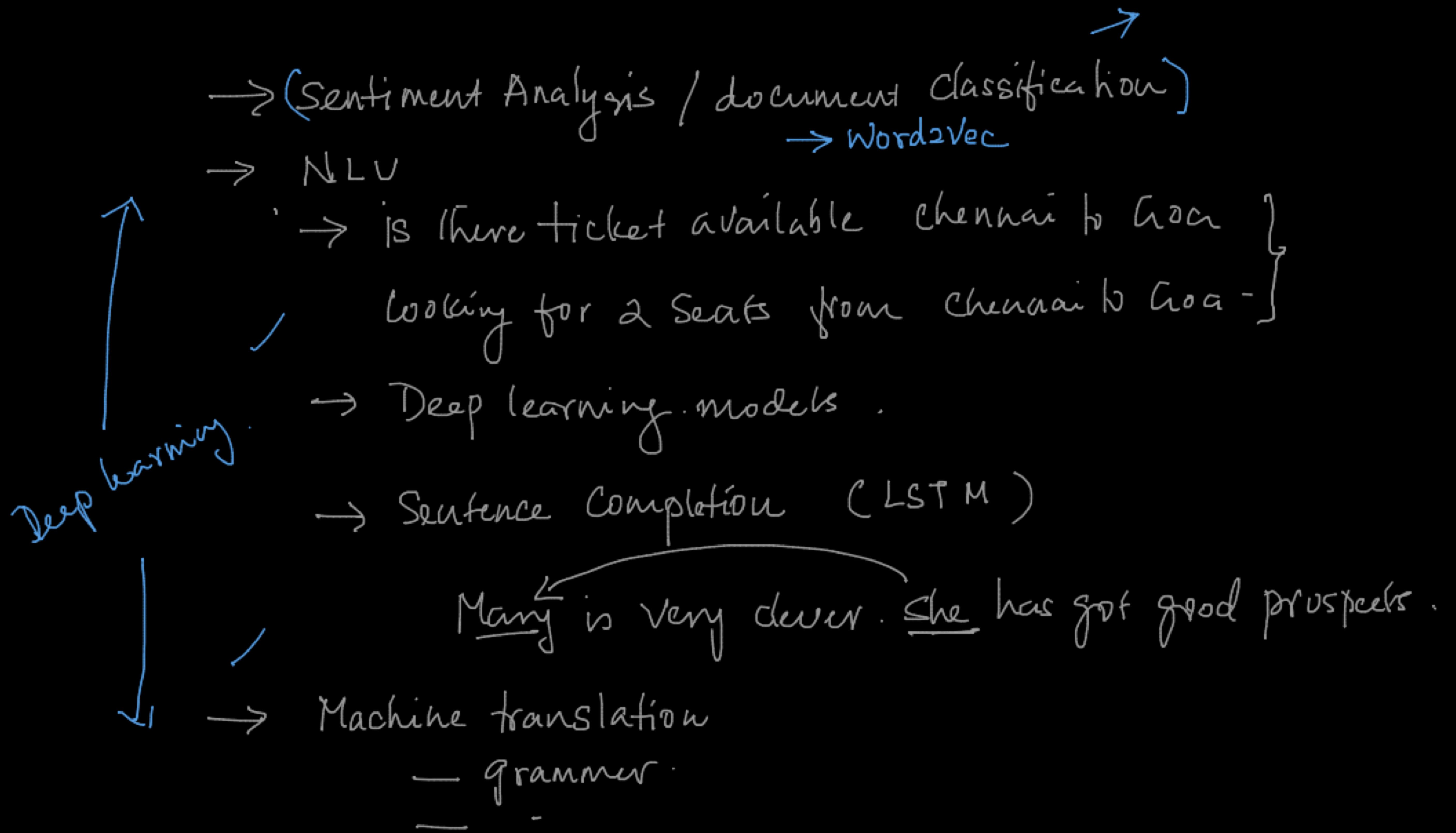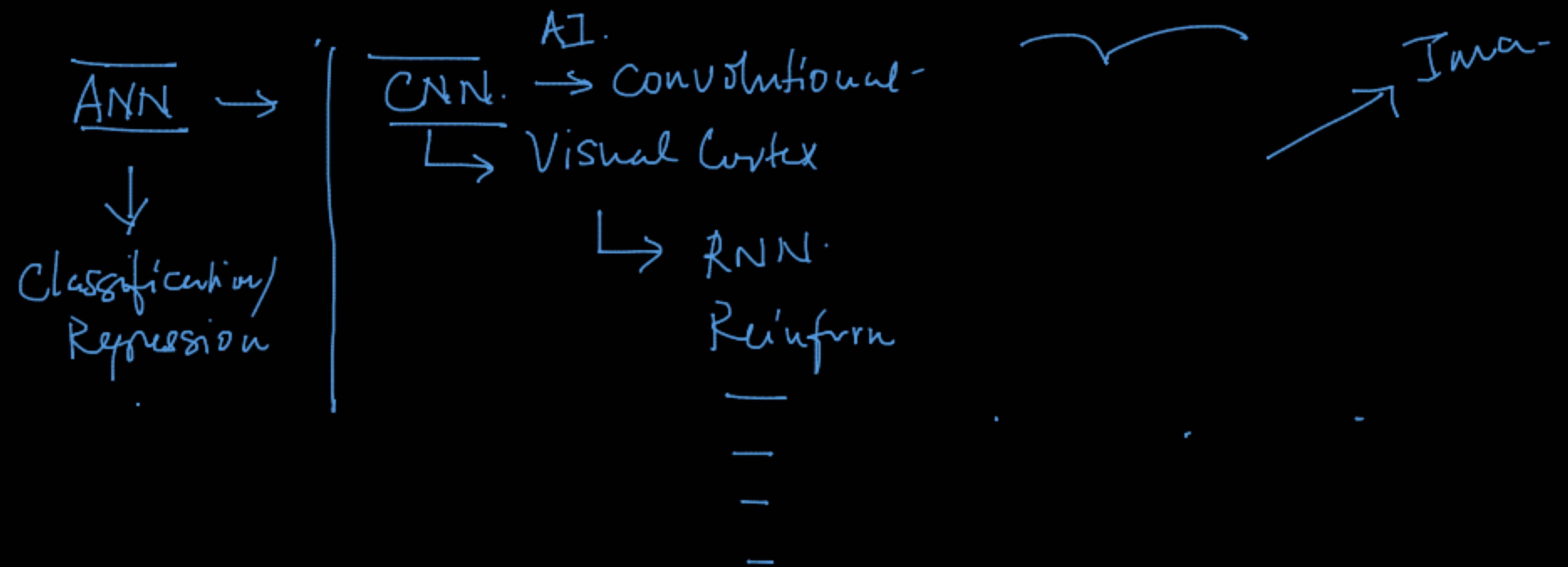
$$\text{IDF} = \frac{1}{\text{DF}}$$

✓ D1: The car is driven on the road. ✓
✓ D2: The truck is driven on the highway

|  | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ | $x_6$ | $x_7$ | $x_8$ | y |
|---|---|---|---|---|---|---|---|---|---|
| → D1 | 0 | 0.043 | 0 | 0 | 0 | 0 | 0.043 | 0 | +ve |
| D2 | 0 | 0 | 0.043 | 0 | 0 | 0 | 0 | 0.043 | -ve |

**TF-IDF**

|  | Term Freq | | DF | IDF | log(IDF) | D1 | D2 |
|---|---|---|---|---|---|---|---|
|  | D1 | D2 |  |  |  |  |  |
| The | 2/7 | 2/7 | 2/2 | 1 | 0 | 0 — | 0 — |
| car | 1/7 | 0 | 1/2 | 2 | 0.3 | 0.043 | 0 — |
| truck | 0 | 1/7 | 1/2 | 2 | 0.3 | 0 | 0.043 |
| is | 1/7 | 1/7 | 2/2 | 1 | 0 | 0 | 0 |
| driven | 1/7 | 1/7 | 2/2 | 1 | 0. | 0 | 0 |
| on | 1/7 | 1/7 | 2/2 | 1 | 0 | 0 | 0 |
| road | 1/7 | 0 | 1/2 | 2 | 0.3 | 0.043 | 0 |
| highway | 0 | 1/7 | 1/2 | 2 | 0.3 | 0 | 0.043 |

$\rightarrow$ (Sentiment Analysis / document Classification)

$\rightarrow$ Word2Vec

$\rightarrow$ NLU

$\quad$ ' $\rightarrow$ is there ticket available chennai to Goa $\Big\}$

$\quad$ looking for 2 seats from chennai to Goa $-\Big\}$

$\rightarrow$ Deep learning models .

$\rightarrow$ Sentence Completion ( LSTM )

Deep learning

Many is very clever. She has got good prospects.

$\rightarrow$ Machine translation

$\quad$ — grammer .

$\quad$ —

$$\overline{ANN} \longrightarrow$$

$$\downarrow$$

Classification/
Regression

.

AI.

$$\overline{CNN.} \longrightarrow Convolutional-$$

$$\hookrightarrow Visual \; Cortex$$

$$\hookrightarrow RNN.$$

Reinforn

—
—
—
—

Ima-

Text processing

$$512 \times 512$$

$$1024 \times 1024 \Rightarrow 10,00\,000$$

$$\Rightarrow 1 \; Million$$

1000

$$\underline{Message} \qquad \underline{Label.}$$

| Message | Label |
|---------|-------|
| How r u ? | Ham |
| Trip for two | Spam |
| ⋮ | ⋮ |

$$Vect = CountVectorizer()$$
$$\Rightarrow Vect.fit\_transform(x_1 - \cdot)$$

$$x_1 \quad x_2 \quad x_3 \quad - - \cdot \quad \cdot \quad x_{7082} \quad | \quad Label.$$

$$clf = linearSVC()$$
$$= clf.fit( \quad )$$
$$= clf.predict( \quad )$$

$$Pipeline \Rightarrow model \cdot (X-tm)$$

Msg | Lub → [ TF-IDF ] — [ SVC. ] → ⋅¹