## Last session:

① Text Preprocessing
② Text Representation
    ① One Hot
    ② BOW
    ③ N-gram

* TFIDF
* word2vec

---

* Agenda:
① TFIDF
② word2vec
③ How to perform text classification
    using Machine Learning
       → Text pre
      → Text Representation

Tommonrow:  RNN

① TFIDF ⟶ Term frequency - Inverse Document Frequency

Drawbacks → ① One Hot — Sparce metria
            ↳ Dimention issue
          ↳ sementic meaning not captured
        ↳ OOV

BOW → Sparce metrix
→ sementic meaning

1, 2, 0

① This movie is Amazing! Amazing Performance by SRK

— — — —
—

most frequent    word will get    more weight    [machine] → 1:h
                                                            (0 heart)
(0.22  0.8)

**TFIDF**

For a term $i$ in document $j$:

$$w_{i,j} = tf_{i,j} \times \log\left(\frac{N}{df_i}\right)$$

$tf_{i,j}$ = number of occurrences of $i$ in $j$
$df_i$ = number of documents containing $i$
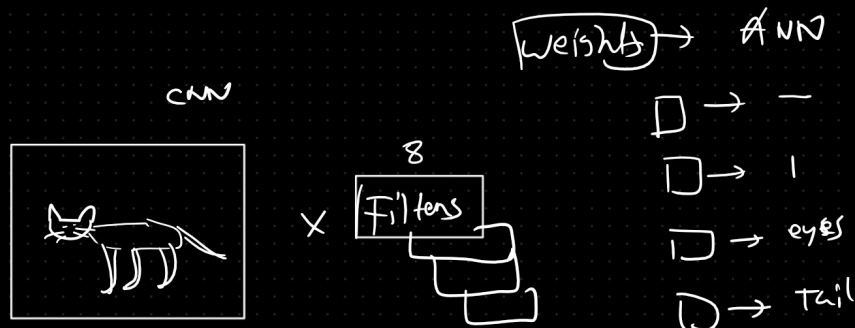$N$ = total number of documents

⟶

| O × O |
| ⇒ O  |

0 ⇒
0.85 × 0.5
⇒ ⌐ ⌐ ⇒  ↗

| This | mail | is | Amaz | perfo | by | | SRK |
|------|------|----|------|-------|----|--|-----|
| 0.85 | 0.75 | 0.5 | (0.99) | 0.7 | 0.4 | | 0.3 |
| — | | — | | — | — | | — |
| | — | — | | | | | |
| | | ↑ | | | | | |

→ postive

→ 100

Ann → |weight| →

② word2vec: It is using some kinds of NN to generate different features from a word

|weight| → ANN

CNN



         8
   × [Filters]

□ → —
□ → 1
□ → eyes
□ → tail

I am the [king] in this Universe

| Features | King | Queen | Man | Women | Monkey |
|----------|------|-------|-----|-------|--------|
| gender | 1 | 0 | 1 | 0 | 1 |
| wealth | 1 | 1 | 0.3 | 0.2 | 0 |
| power | 1 | 0.7 | 0.5 | 0.2 | 0 |
| weight | 0.8 | 0.5 | 0.7 | 0.5 | 0.3 |
| speak | 1 | 1 | 1 | 1 | 0 |

→ Word2vec

NM →

BP

king → $[1, 1, 1, 0.8, 1]$ —

Queen ⇒ $[0, 1, 0.7, 0.5, 1]$ —

5D

substract



| Features | King | Queen | Man | Women | Monkey |
|----------|------|-------|-----|-------|--------|
| gender | 1 | 0 | 1 | 0 | 1 |
| wealth | 1 | 1 | 0.3 | 0.2 | 0 |
| power | 1 | 0.7 | 0.5 | 0.2 | 0 |
| weight | 0.8 | 0.5 | 0.7 | 0.5 | 0.3 |
| speak | 1 | 1 | 1 | 1 | 0 |

Word2vec



my name is Bapy I am a Batal

A Game Of Thrones Book One of A Song
of Ice and Fire By George R. R. Martin
PROLOGUE "We should start back," Gared
urged as the woods began to grow dark
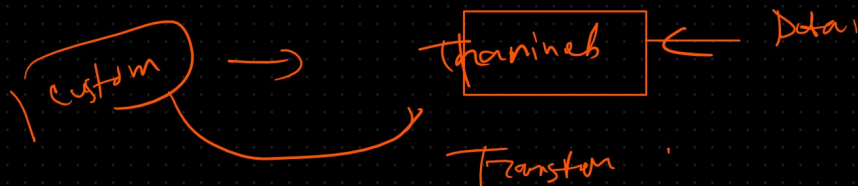around them. "The wildlings are dead."

(Rdv) 8 G 02

Simple pre procer sensing

n  0    A Game Of Thrones Book One of A Song

S  0    of Ice and Fire By George R. R. Martin

l  ③    PROLOGUE "We should start back," Gared

[ "A", "Games"

Jensim → word2 vee

custom → thanineb ← Detai

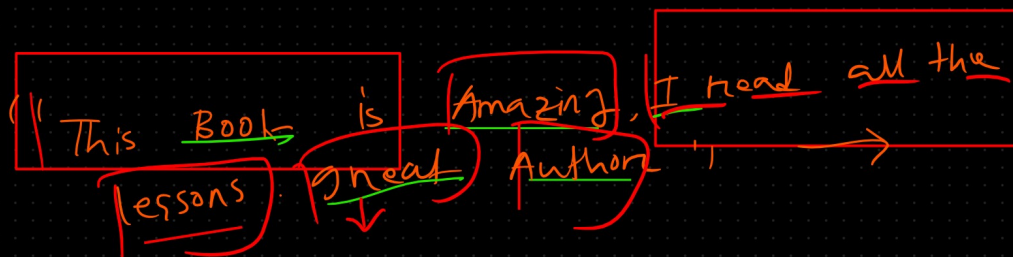Transform

word2 vect

100 →    | 100,    King    100 D

A Game Of Thrones Book One of A Song
of Ice and Fire By George R. R. Martin
PROLOGUE "We should start back," Gared
urged as the woods began to grow dark
around them. "The wildlings are dead."

→ 100

Amazin, 100

c X This -2

"This Book is Amazing, I read all the

lessons great Author

frequent Accuracy comel

con