# Day 3 - NLP

## Agenda

1. BOW {Bag of words}
2. Tf Idf
3. Practical Implementation
4. Quiz

## Spam Classifier

### Word2vec

1. Embedding Layer
2. Word2vec → CBOW ✓ → Skip gram ✓
3. Architecture
4. Practical Problem
5. Glove

Bag of words → Text → Vectors

---

① 

Sent 1 → (He) (is) (a) g(Good↑)ood boy

Sent 2 → (She) (is) (a) good girl

Sent 3 → Boy (and) girl (are) good

② Stopwords ⟹ lowering

| | | |
|---|---|---|
| Sent1 | good | boy |
| Sent 2 | good | girl | cat |
| Sent 3 | boy | girl | good |

---

③ **Frequency (Vocabulary)**

| | frequency |
|---|---|
| good | 3 |
| boy | 2 |
| girl | 2 |

Not all Similar

④

| | $\underset{f_1}{good}$ ⟺ | $\underset{f_2}{boy}$ | $\underset{f_3}{girl}$ |
|---|---|---|---|
| Sent 1 | ① ⟷ | ① | 0 |
| Sent 2 | 1 ✓ | 0 | 1 ✓ |
| Sent 3 | 1 | 1 | 1 |

---

opposite {
The food is good
The food is not good
}

| food | good | ✗ | not | ~~the~~ |
|---|---|---|---|---|
| 1 | 1 | 1 | 0 | 1 |
| 1 | 1 | 1 | 1 | 1 |

⟹ { TF - IDF }

⟹ Similar

# Term Frequency — Inverse Document Frequency        TF - IDF

Sent 1 :  good  boy ✓

Sent 2 :  good  girl ✓

Sent 3 :  boy girl good ✓

$$\text{Term Frequency} = \frac{\text{No. of rep of words in sentence}}{\text{No. of words in sentence}}$$

↓ Sentences

X

$$IDF = \log_e \left( \frac{\text{No. of sentences}}{\text{No. of sentences containing the word}} \right)$$

↓ words

## Term Frequency          *          ## Inverse Document Frequency

|  | Sent 1 | Sent 2 | Sent 3 |  | Words | IDF |
|------|--------|--------|--------|---|-------|-----|
| good | $\frac{1}{2}$ | $\frac{1}{2}$ | $\frac{1}{3}$ | * | good | $\log_e(3/3) = 0$ |
| boy  | $\frac{1}{2}$ | $0$ | $\frac{1}{3}$ |  | boy | $\log_e(3/2)$ |
| girl | $0$ | $\frac{1}{2}$ | $\frac{1}{3}$ |  | girl | $\log_e(3/2)$ |

⇓

O/P

|  | f1 good | f2 boy | f3 girl |
|--------|---------|--------|---------|
| Sent 1 | $0$ ✓ | $\frac{1}{2} \times \log_e(3/2)$ ✓ | $0$ |
| Sent 2 | $0$ ✓ | $0$ | $\frac{1}{2} \times \log_e(3/2)$ |
| Sent 3 | $0$ ✓ | $\frac{1}{3}(\log_e(3/2))$ ✓ | $\frac{1}{3}\log_e(3/2)$ ✓ |

## Advantages

① Intuitive

② Word Importance is getting capture }

## Disadvantage

① Sparsity

② Out of vocabulary

good ✓          10                    max.f = 3              good    boy    girls

boy ✓           7                                         [                    ]

girls ✓         6                    Assignment

try             5                    Take any Text Dataxt

|               1                    ┌ Apply BOW, TF-IDF, ┐ ✓
.               1                    │                    │ ✓
.               1                    └    ngram           ┘
.               1
.               2
<3

                ┌─────────────────────────────┐
                │ KrishnaiK06@gmail.com       │
                └─────────────────────────────┘