# TF-IDF [Intuition]

→ Term frequency & inverse document frequency.

Term frequency (TF) = $\dfrac{\text{No. of repetitⁿ of words in sentences}}{\text{No of words in sentence.}}$

Inverse Document frequency (IDF) = $\log\left(\dfrac{\text{No. of sentences}}{\text{No. of sentences containing words}}\right)$

\#. finally $\boxed{TF \times IDF}$ is used to generate vector

| | | | word | frequ |
|---|---|---|---|---|
| $S_1$ | good boy | | good | 3 |
| $S_2$ | good girl | | boy | 2 |
| $S_3$ | good boy girl. | | girl | 2 |

**TF**

| | $S_1$ | $S_2$ | $S_3$ |
|---|---|---|---|
| good | 1/2 | 1/2 | 1/3 |
| boy | 1/2 | 0 | 1/3 |
| girl | 0 | 1/2 | 1/3 |

vectors →

~~good~~

TF-IDF →

**IDF**

| words | IDF |
|---|---|
| good | $\log(3/3) = 0$ |
| boy | $\log(3/2)$ |
| girl | $\log(3/2)$ |

| | $f_1$ | $f_2$ | $f_3$ | op |
|---|---|---|---|---|
| | good | boy | girl | |
| $S_1$ | 0 | $\frac{1}{2}\log(\frac{3}{2})$ | 0 | |
| $S_2$ | 0 | 0 | $\frac{1}{2}\log(\frac{3}{2})$ | |
| $S_3$ | 0 | $\frac{1}{3}\log(\frac{3}{2})$ | $\frac{1}{3}\log(3/2)$ | |

used to train my also as it is dependent feature.

# TF-IDF

```python
import nltk
Pagraph = _____
                _____
                ___

# cleaning
import re
from nltk.stem import PorterStemmer
from nltk.stem import WordNetLammatizer
from nltk.corpus import stopwords


stemmer = PorterStemmer()
Lammatizer = WordNetLemmatizer()
sentences = nltk.sent_tokenize(paragraph)
corpus = []


for i in range(len(sentences)):
    review = re.sub('[^a-zA-Z]', ' ', sentences[i]))
    review = review.lower()
    review = review.split()
    review = [stemmer.stem(word) for
        word in review if word not in set(stopwords.
                words('english'))]
    review = ' '.join(review).
    corpus.append(review)


# Calculate the TF-IDF
from sklearn.feature_extraction.text import TfidfVectorizer.
cv = TfidfVectorizer()
x = cv.fit_transform(corpus).toarray().
```