## Agenda.

NLP- Natural Language Processing

Non-numeric
Features

| Categorical | Ordinal | Pure Text |
|---|---|---|
| Label Encoding | (Replace values | data. |
| OneHotencoding | with relevant | |
| | numeric value) | |

The scope of NLP is to handle pure txt data.

Text → [NLP] → Numeric
Feature        feature
                          ↓
                     Classifi^n
                     algs
                          ↓
                     model

# NLP program rules

① No punctuations
② No stopwords

**Python Basics**

Preprocessing ← NHK package

Text Feature →
| ① Remove punctuations
| ② Convert sentences into words
| ③ Remove Stopwords.

↓

Create BagOfWords }

Create a vocab. dictionary which contains text and its count
eg. (Hello, 27223) ←

↓

Convert count into frequency both (0-1)

TF - IDF }

Term Frequency - Inverse Document Frequency (Hello, 0.8)

↓

Create feature array.

↓

Label ⟶ ml algo ⟶ model

**Naive Bayes**

*(vertical text, left side:)* SMS Spam Classifier

*(vertical text:)* SKLEARN

"coin
lottery
guaranteed" $\longrightarrow$ [ model ] $\longrightarrow$ Spam

eg:

dataset = [ "welcome to Simplilearn",

"welcome welcome welcome",

" Simplilearn, Simplilearn welcome"]

$\downarrow$
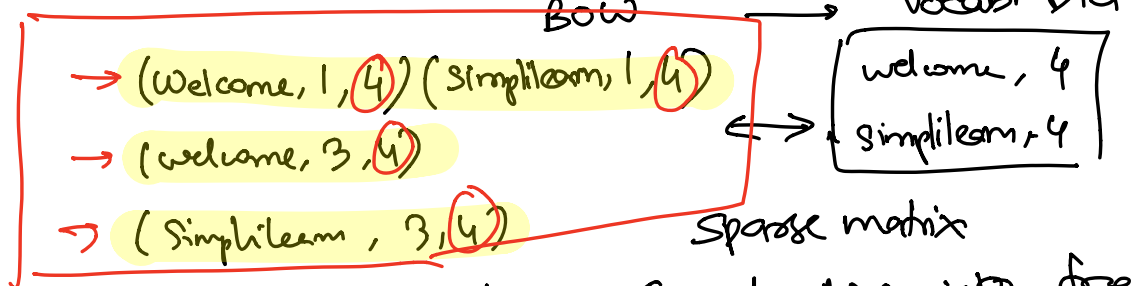
① Reprocess.

(welcome, Simplilearn)

(welcome, welcome, welcome)

(Simplilearn, Simplilearn, Simplilearn)

$\downarrow$

BOW $\longrightarrow$ Vocab. Dict

→ (Welcome, 1, 4) (Simplilearn, 1, 4)

→ (welcome, 3, 4)

→ (Simplilearn, 3, 4)

| welcome, 4 |
| Simplilearn, 4 |

Sparse matrix

Convert BOW into freq
using TFIDF

$\downarrow$

TFIDF Transformer

→ fit(bow) →

{ (welcome (1/4)) (Simplilearn, 1/4)
(welcome, 3/4)
(Simplilearn, 3/4) sparse matrix

feature

| 0.25, 0.25 |
| 0.75 |
| 0.75 |

$\longleftarrow$

sparse matrix