| Notes

text of the tweet
Date ......→ polarity
extract (X) and (Y) label

raw_csv_file →

Data Preprocessing →
→ lower_case
→ remove VRL's
→ remove mentions
→ keep only letters
→ double white spaces or new lines remove

df ["clean_text"]

Tokenize

df ["tokens"]

I, am, and, is, or
remove common words (stopwords) → pronoun (X)
→ verbs (X)
removed

df ["tokens_nostop"]

Tweet → fixed length numeric vector

TF-IDF

$X = \begin{bmatrix} 0 & 1 & \cdots \\ \vdots & & \cdots \end{bmatrix}$  sparsh matrix  (TF-IDF matrix)

Dataset → 80% Train → X_train, Y_train →
→ 20% Test → X_test  Y_test

logistic regression

Training model  model.fit (X_train, Y_train) ←  ? P X / OR

→ & y pred.

# Prediction on Test dataset
y_pred = model.predict (X_test)
# Accuracy
acc = accuracy_score (Y_test, Y_pred)

0.78

given predicted by model

| Notes ✓

## TF-IDF (Term frequency - Inverse Document frequency)

⤷ converts ~~text~~ sentence ⟶ ~~numbs~~ vector of number while preserving meaning.

**(st) Vocabulary creation**

⤷ scans all tweets and builds dictionary of words.

love → 0
hate → 1
machine → 2
learning → 3
⋮

TF-IDF shape $(1,048,572, 15000)$

↑ Tweets    ↑ unique words.

⟶ importance

$$TF(word, Tweet) = \dfrac{count\ of\ word\ in\ a\ tweet \uparrow}{total\ words\ in\ a\ tweet}$$

$$IDF\left(\begin{array}{c}Inverse\ Document\\frequency\end{array}\right) = log\left(\dfrac{total\ \overset{tweets}{\cancel{documents}}}{tweets\_containing\_word \uparrow}\right)$$

⟶ How-rare      IDF↑ ⟶ rare words
                TF↑ ⟶ important

TF×IDF

(i) love machine learning $= [0.32, 0.52, 0.52, 0.11]$

max-features $= 5500$ (most frequent words)

$=$