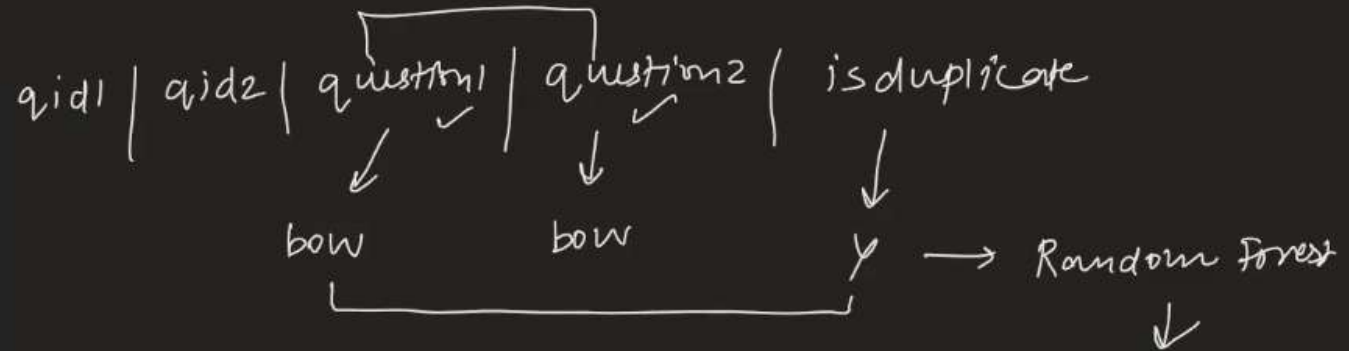


Friday, February 11, 2022 7:05 AM



- 1 q1 len → char length of q1
- 2 q2 len → char length of q2
- 3 q1 words → # words in q1
- 4 q2 words → # words in q2
- 5 words common → # of common unique words
- 6 words total → Total # words in q1 + Total # of words in q2
- 7 word share → word common / word total

Token Features

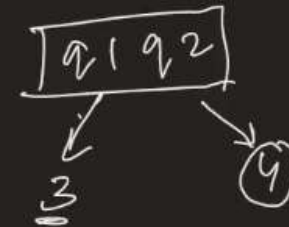
- 1) cwc - min -
- 2) cwc - max -
- 3) csc - min -
- 4) csc - max -
- 5) ctc - min -
- 6) ctc - max -
- 7) last_word - eq
- 8) first_word - eq

Length Based features

tokens, words, stop words is of
 ↓ ↓ ↓
 5 3 2
 [sachin] [is] [a] [great] [batsman]

$$\rightarrow \frac{\# \text{ common words}}{\min(\text{words}(q_1, q_2))}$$

(3) ← min max



q_1 q_2
 1 1
 5 6

$$\rightarrow \frac{\# \text{ common stop words}}{\min(\text{stop words}(q_1, q_2))} \leftarrow 5$$

10 ← min max

q_1 q_2
 1 1
 10 15

$$\rightarrow \frac{\# \text{ common tokens}}{\min(\text{tokens}(q_1, q_2))} \leftarrow 10$$

10 ← min max

predict

	0	1
actual 0	✓	✗
1	✗	✓

$AD \rightarrow 0$
 $PD \rightarrow 1$

$\left\{ \begin{array}{l} \text{random} \\ \text{xgbost} \end{array} \right.$



1) Increase data (42) → 30K
↓
RAM

→ Increase RAM
→ Cloud platform
→ incremental learning
↓
Vortex / Dask

2) Preprocessing → Stemming

3) Apply more algorithms

→ SVM Logistic,
→ Hyperparameters tuning
→ cross validation

4) More feature (22) →

5) Bag of word

→ tfidf
→ word2vec
→ tfidf weighted w2v

6) Deep learning → 4D
↓
data loaders

