

DESCRIPTION OF FEATURES :

Based on trained data 5 Features has been used to initialize.

→ Features have been chosen based on frequency, i.e. most common occurring words

i) Common words in english like "which", "if", "is", "at", "then" ... etc

ii) Dutch alphabets that are unique  
ä, é, ï, ò, ù

iii) Dutch numbers like "een", "twee", ....

iv) English numbers like "one", "two", ....

v) Dutch common words "goed", "zif", ....

## DECISION TREE !

- i) Four parameters have been used level, maximum level, already visited nodes and output file
- ii) Level with one has been treated as root node. Visited node has been initialized. Using the total count of features a matrix has been created using which entropy of each column has been calculated.
- iii) A decision tree will be created based on the entropy values. A visited set is also maintained to avoid redundancy.
- iv) Maximum level of 5 has been used as operation beyond level 5 has not produced anything different.

v) Decision tree outputs have been included in the examples folder within the folder Submitted.

### ADABOOST :-

i) Condition being consider as :

Weight  $\propto \frac{1}{\text{length of data frame}}$

Weight is inversely proportional to length of data frame.

ii) Here 10 decision tree have been created, trained and based on that it guess either true or false which leads us to the actual answer.

iii) Weight values will be updated based on the error value calculated.

iv) If the decision tree gives the expected output, the updated weight will be less than the actual weight, whereas the updated weight will be more for wrong decisions.

v.) This process keeps operating for all the decision trees. Error rate's will be updated based on the training of every decision tree, and based on the error rates adaboost gives the output through the comparison of values that is english or Dutch. I have all the decision trees conclusively provide adaboost result.