

## Advanced Features

### 1. Token Features

- **cwc\_min**: This is the ratio of the number of common words to the length of the smaller question
- **cwc\_max**: This is the ratio of the number of common words to the length of the larger question
- **csc\_min**: This is the ratio of the number of common stop words to the smaller stop word count among the two questions
- **csc\_max**: This is the ratio of the number of common stop words to the larger stop word count among the two questions
- **ctc\_min**: This is the ratio of the number of common tokens to the smaller token count among the two questions
- **ctc\_max**: This is the ratio of the number of common tokens to the larger token count among the two questions
- **last\_word\_eq**: 1 if the last word in the two questions is same, 0 otherwise
- **first\_word\_eq**: 1 if the first word in the two questions is same, 0 otherwise

### 2. Length Based Features

- **mean\_len**: Mean of the length of the two questions (number of words)
- **abs\_len\_diff**: Absolute difference between the length of the two questions (number of words)
- **longest\_substr\_ratio**: Ratio of the length of the longest substring among the two questions to the length of the smaller question

### 3. Fuzzy Features

- **fuzz\_ratio**: fuzz\_ratio score from fuzzywuzzy
- **fuzz\_partial\_ratio**: fuzz\_partial\_ratio from fuzzywuzzy
- **token\_sort\_ratio**: token\_sort\_ratio from fuzzywuzzy
- **token\_set\_ratio**: token\_set\_ratio from fuzzywuzzy

[FuzzyWuzzy: Fuzzy String Matching in Python - ChairNerd](#)

<https://stackoverflow.com/a/19794953> - Search

[https://en.wikipedia.org/wiki/Wikipedia%3aList\\_of\\_English\\_contractions](https://en.wikipedia.org/wiki/Wikipedia%3aList_of_English_contractions)