## QUESTION - 6

| | UNIGRAM-1000 | UNIGRAM-1000 | BPE-1000 | BPE-2000 | mBERT - 1000 | mBERT - 2000 | Indic-BERT -1000 | Indic-BERT -2000 | WHITE SPACE |
|---|---|---|---|---|---|---|---|---|---|
| **PRECISION** | 0.01724 1379310 344827 | 0.022988 50574712 6436 | 0.01149 4252873 563218 | 0.0172 413793 103448 27 | 0.0057 471264 367816 09 | 0.005 74712 64367 81609 | 0.01149 4252873 563218 | 0.01149 4252873 563218 | 0.04022 9885057 471264 |
| **RECALL** | 0.00384 1229193 3418692 | 0.005847 95321637 4269 | 0.00247 2187886 279357 | 0.0042 134831 460674 16 | 0.0012 406947 890818 859 | 0.001 24069 47890 81885 9 | 0.00318 9792663 476874 | 0.00318 9792663 476874 | 0.01548 6725663 716814 |
| **F-SCORE** | 0.00628 2722513 0890054 | 0.009324 00932400 9324 | 0.00406 9175991 861648 | 0.0067 720090 293453 73 | 0.0020 408163 265306 124 | 0.002 04081 63265 30612 4 | 0.00499 3757802 746566 | 0.00499 3757802 746566 | 0.02236 4217252 396165 |

**COMPARISON** :

The PrecisIon, Recall and F-score value is highest for White space tokenizer because other tokenizers are giving tokens that are even dividing the single words. That is causing a large mismatch with the word tokens that I have marked in question 3 as per definition of word group. And white space tokenizer is words from spaces which increases its match from ground truth list. Hence, the true positive high for white space tokenizer.

But according to the description of tokenizers we should get the highest precision for Indic-BERT tokenizer as it is trained for Hindi language.