**Table: Justification of 11 features selection to predict multiple diseases (continue).**

| Justification for selection 11 features | | | |
|---|---|---|---|
| **Feature Category** | **Feature** | **Feature Score** | **Reasons** |
| Lexical | Noun Count | Correlation Coefficient: 0.157469<br>Chi-Square: 106.006848<br>Mutual Information: 0.156650 | Noun count is strongly related to text meaning, as nouns often serve as key elements in sentence structure. The chi-square value of 106 and the mutual information score of 0.156650 confirm its significant role in distinguishing the target variable, making it a valuable feature for classification. |
| Semantic | Polarity | Correlation Coefficient: 0.150401<br>Chi-Square: 118.543829<br>Mutual Information: 0.763562 | Polarity reflects sentiment, which is essential for understanding text's emotional tone. The strong chi-square result and high mutual information (0.763562) emphasize its importance in predicting the target variable, supporting its inclusion as a top feature. |
| Syntactic | Fragments Proportion Score | Correlation Coefficient: 0.129982<br>Chi-Square: 33.211240<br>Mutual Information: 0.040863 | Justification: The proportion of fragments indicates the complexity and structure of text. While the correlation is modest, the chi-square value (33.211240) and mutual information (0.040863) show that fragments significantly differentiate the target variable, thus validating its importance. |
| Lexical | Long Words Count | Correlation Coefficient: 0.069164<br>Chi-Square: 72.522411<br>Mutual Information: 0.098928 | Longer words are associated with more complex text, which can be indicative of certain target classes. The chi-square and mutual information results support its relevance, even though the correlation coefficient is relatively lower, suggesting its importance in the classification task. |
| Semantic | Irony Sarcasm Count | Correlation Coefficient: 0.060460<br>Chi-Square: 7.413642<br>Mutual Information: 0.025005 | Irony and sarcasm often provide nuanced information about text meaning. Although the correlation is low, the chi-square result (7.413642) and mutual information (0.025005) demonstrate its utility in distinguishing the target variable, making it relevant for the feature set. |
| Lexical | Average Word Length | Correlation Coefficient: 0.055547<br>Chi-Square: 57.588631<br>Mutual Information: 0.998094 | Justification: Average word length can indicate the complexity of the text. The high mutual information value (0.998094) and reasonable chi-square (57.588631) support its contribution to predicting the target variable, ensuring its place in the top features. |
| Readability | Automated Readability Index | Correlation Coefficient: 0.048720<br>Chi-Square: 1359.460186<br>Mutual Information: 1.238450 | The Automated Readability Index measures text readability, which is a key determinant of text classification. The very high chi-square value (1359.460186) and mutual information (1.238450) confirm its strong predictive power, supporting its inclusion in the top features. |

**Table: Justification of 11 features selection to predict multiple diseases.**

| Justification for selection 11 features (continue) | | | |
|---|---|---|---|
| **Feature Category** | **Feature** | **Feature Score** | **Reasons** |
| Syntactic | Syntax Tree Depth Variability Score | Correlation Coefficient: 0.036676<br>Chi-Square: 32.145670<br>Mutual Information: 0.734970 | The depth and variability of a syntax tree reflect text structure and complexity. While the correlation is modest, the chi-square (32.145670) and mutual information (0.734970) show that this feature significantly contributes to distinguishing the target variable. |
| Semantic | Coreference Resolution Density | Correlation Coefficient: 0.020171<br>Chi-Square: 10.912349<br>Mutual Information: 0.213544 | Justification: Coreference resolution is crucial for understanding text coherence. Despite a low correlation coefficient, the chi-square (10.912349) and mutual information (0.213544) values highlight its importance in determining text meaning and differentiating target categories. |
| Semantic | Relationships Variation | Correlation Coefficient: 0.015332<br>Chi-Square: 1.834900<br>Mutual Information: 0.023761 | Variation in relationships provides insight into text structure and meaning. While the correlation coefficient is weak, the chi-square and mutual information results confirm that this feature contributes to distinguishing the target variable. |
| Lexical | Lexical Diversity Score | Correlation Coefficient: 0.005532<br>Chi-Square: 9.123765<br>Mutual Information: 0.106017 | Lexical diversity measures the variety of vocabulary in a text. Although the correlation is low, its chi-square value (9.123765) and mutual information (0.106017) suggest that it helps differentiate between text classes, making it relevant for the analysis. |
| 2 Features not considered (despite positive correlation coefficient score) | | | |
| Syntactic | Clause Boundary Standard Deviation | Correlation Coefficient: 0.001656<br>Chi-Square: 0.822956<br>Mutual Information: 0.026603 | While it has a positive correlation with the target variable, both the chi-square (0.822956) and mutual information (0.026603) are very low, indicating that this feature does not significantly contribute to distinguishing the target. As a result, it was removed from the top features. |
| Lexical | Lexical Sophistication Score | Correlation Coefficient: 0.000953<br>Chi-Square: 0.210168<br>Mutual Information: 0.12842 | Justification: This feature also shows positive correlation with the target, but the low chi-square (0.210168) and mutual information (0.12842) values indicate that it does not have a substantial impact on the target variable classification. Thus, it was excluded from the selected features. |