**Sanjivani Rural Education Society's**
**Sanjivani College of Engineering, Kopargaon-423 603**
*(An Autonomous Institute, Affiliated to Savitribai Phule Pune University, Pune)*
*NACC 'A' Grade Accredited, ISO 9001:2015 Certified*

## Department of Computer Engineering
*(NBA Accredited)*

# Subject- Data Mining and Warehousing Lab (CO319)
# Assignment No.5

Prof. S.A.Shivarkar
Assistant Professor
Contact No.8275032712
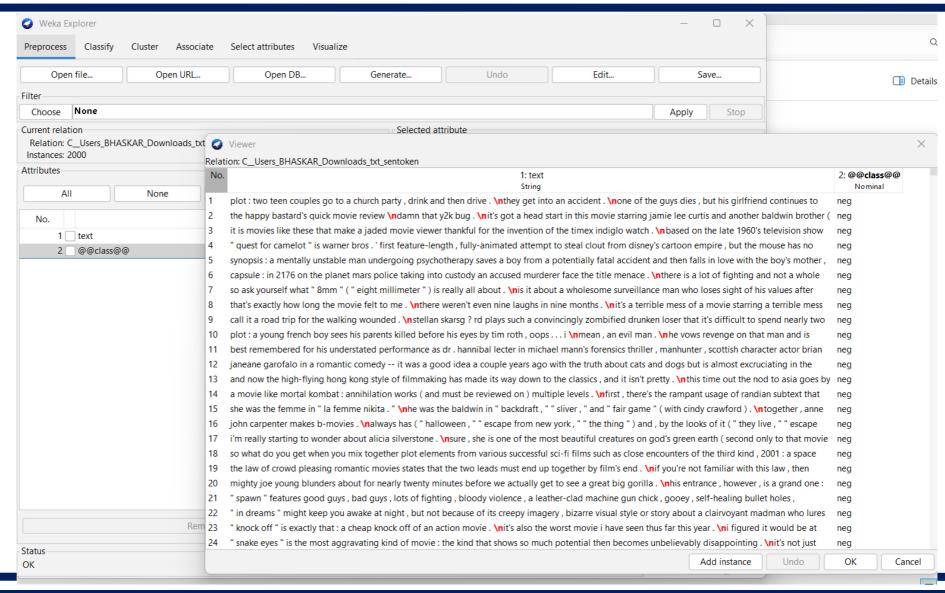Email- shivarkarsandipcomp@sanjivani.org.in

# Problem Statement

- Consider a suitable text dataset. Remove stop words, apply stemming and feature selection techniques to represent documents as vectors. Classify documents and evaluate precision, recall (For Ex. Movie Review Dataset)
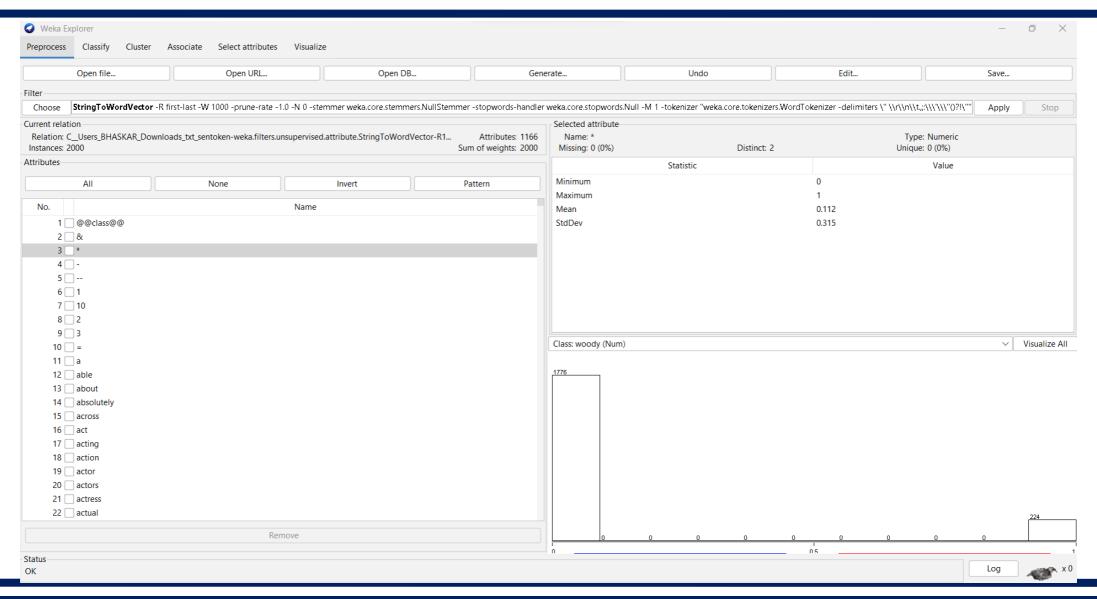
# View the data

Select Class→
Click on
Edit→ View
the data

# Converting string to Numeric attribute
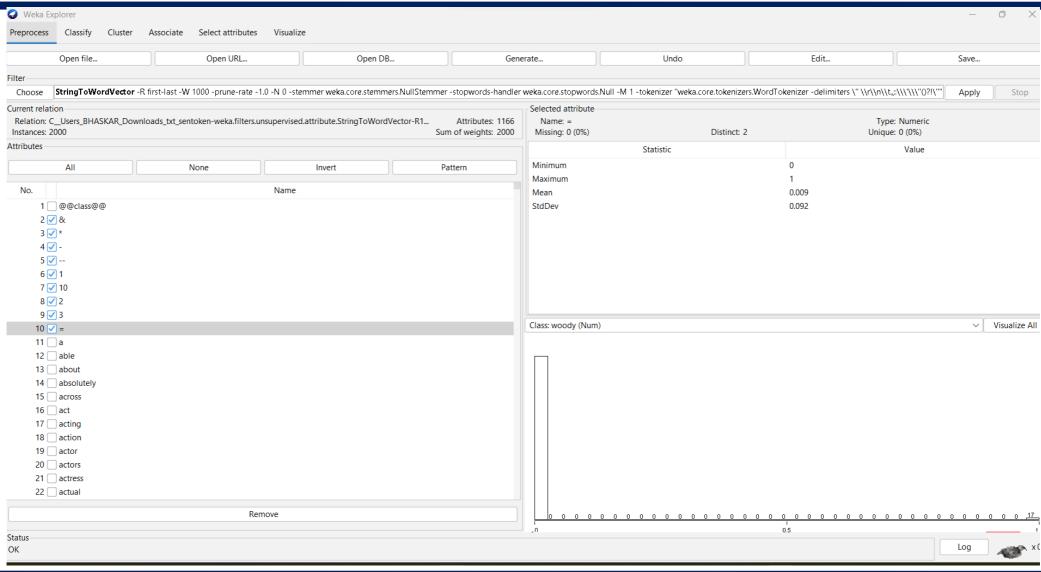
- Select Class→ Choose Filter→ Filter→ Unsupervised→ Attribute→ StrinToWordVector→Apply
- This will convert each word in string field to numeric attribute.
- Name of the attribute is word itself.
- The value of each attribute is 0 (absent) or 1 (present) in the current document.
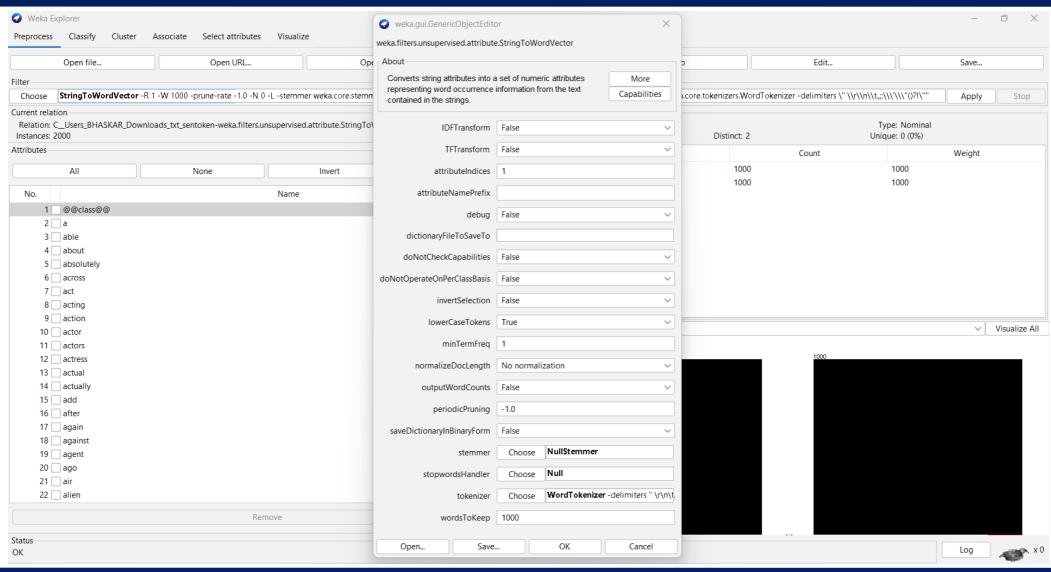
# Converting string to Numeric attribute

# Remove the selected attribute

# Converting string to Numeric attribute

- Click on StringToWordVector in filter.
- Set the parameters attributeIndices=1 and lowerCaseTokens=True, rest keep default.
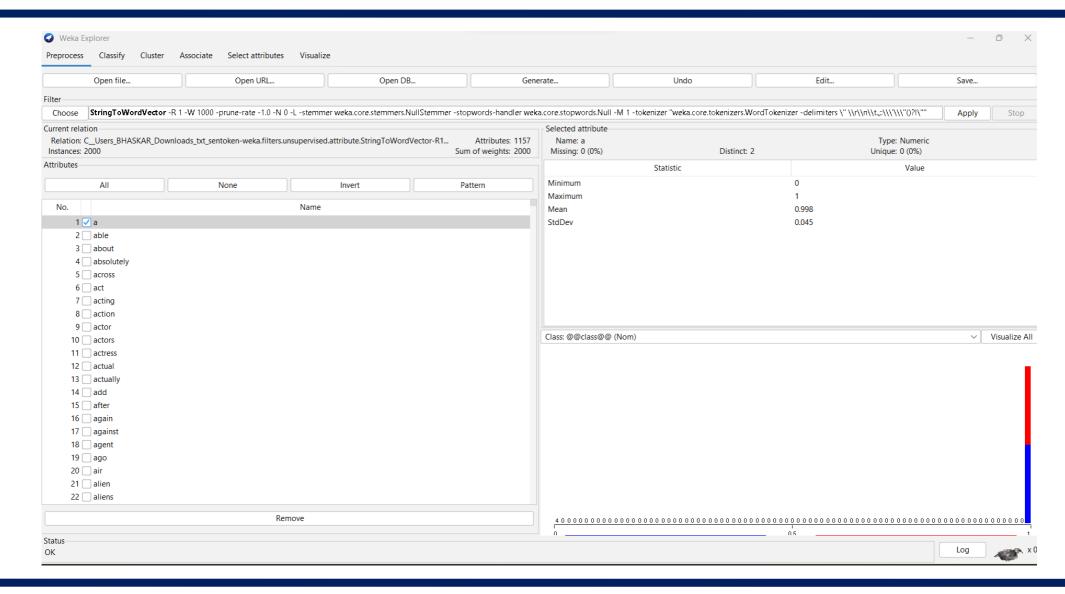- Click on Ok

# Converting string to Numeric attribute

# Move class attribute to the end of WordVector
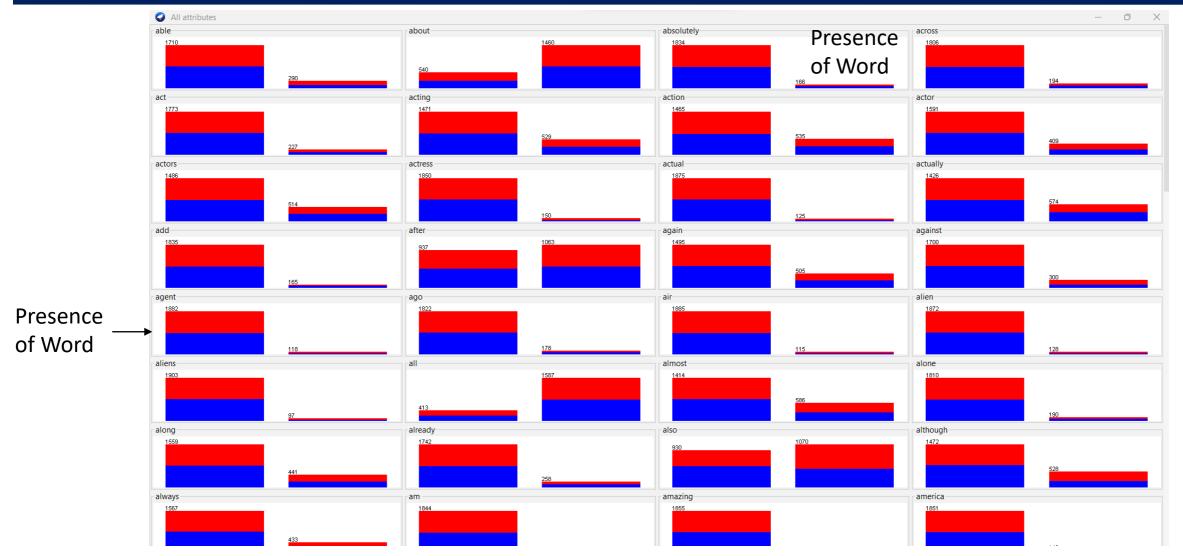
Click on
Edit→ Right
click on
@@class@@
attribute→
Select
attribute as
class

# Remove junk words

- Choose Filter→ Filter→ Unsupervised→ Attribute→ Numeric→Apply
- This will convert each word in string field to numeric attribute.
- Name of the attribute is word itself.
- The value of each attribute is 0 (absent) or 1 (present) in the current document.

# Visualize each attribute as Histogram

# Classification using Naive Bayes Classifier

- Classify→Choose Classifier→Bayes → NaiveBayes→ Select Default option Percentage Split 66% → Start

```
Classifier output
 0                    970.0  965.0
 1                     32.0   37.0
 [total]             1002.0 1002.0



Time taken to build model: 0.24 seconds

=== Evaluation on test split ===

Time taken to test model on test split: 0.44 seconds

=== Summary ===

Correctly Classified Instances         524              77.0588 %
Incorrectly Classified Instances       156              22.9412 %
Kappa statistic                          0.5413
Mean absolute error                      0.2306
Root mean squared error                  0.4363
Relative absolute error                 46.1105 %
Root relative squared error             87.256  %
Total Number of Instances              680

=== Detailed Accuracy By Class ===

              TP Rate  FP Rate  Precision  Recall  F-Measure  MCC     ROC Area  PRC Area  Class
              0.805    0.264    0.752      0.805   0.778      0.543   0.856     0.851     neg
              0.736    0.195    0.792      0.736   0.763      0.543   0.856     0.852     pos
Weighted Avg. 0.771    0.229    0.772      0.771   0.770      0.543   0.856     0.852

=== Confusion Matrix ===

   a    b   <-- classified as
 273   66 |   a = neg
  90  251 |   b = pos
```

# Classification using Naive Bayes Classifier with all attributes

- Classify→Choose Classifier→Bayes → NaiveBayes→ Select Default option Percentage Split 66% → Start

Accuracy 77.0588 % for all attributes

```
Classifier output

  0                    970.0  965.0
  1                     32.0   37.0
  [total]             1002.0 1002.0



Time taken to build model: 0.24 seconds

=== Evaluation on test split ===

Time taken to test model on test split: 0.44 seconds

=== Summary ===

Correctly Classified Instances        524              77.0588 %
Incorrectly Classified Instances      156              22.9412 %
Kappa statistic                         0.5413
Mean absolute error                     0.2306
Root mean squared error                 0.4363
Relative absolute error                46.1105 %
Root relative squared error            87.256  %
Total Number of Instances             680

=== Detailed Accuracy By Class ===

              TP Rate  FP Rate  Precision  Recall  F-Measure  MCC    ROC Area  PRC Area  Class
              0.805    0.264    0.752      0.805   0.778      0.543  0.856     0.851     neg
              0.736    0.195    0.792      0.736   0.763      0.543  0.856     0.852     pos
Weighted Avg. 0.771    0.229    0.772      0.771   0.770      0.543  0.856     0.852

=== Confusion Matrix ===

   a    b   <-- classified as
 273   66 |   a = neg
  90  251 |   b = pos
```
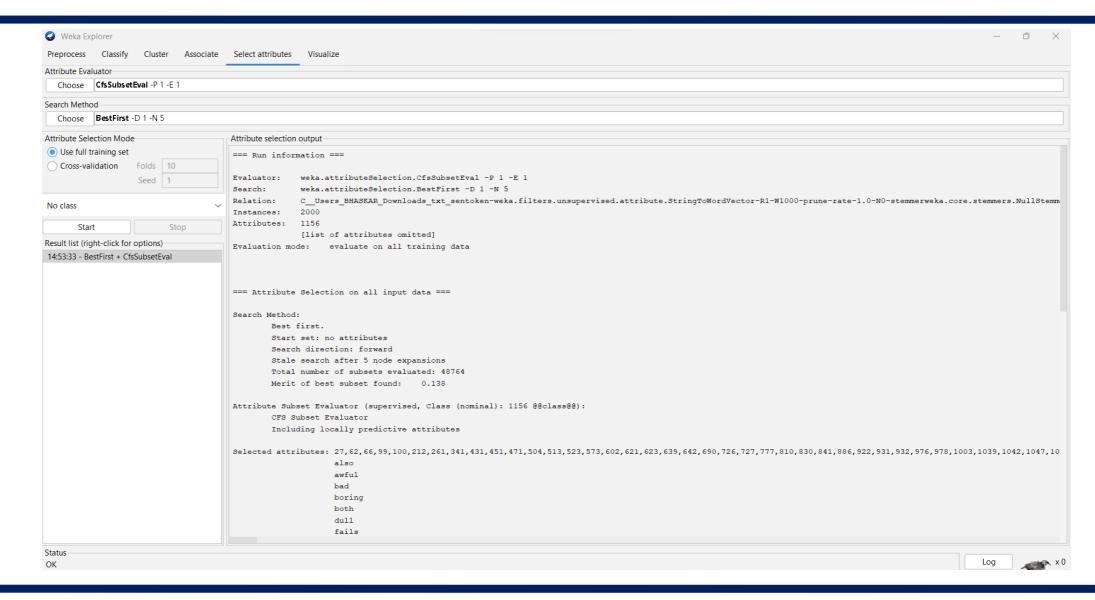
# Attribute Selection

- Select Attributes→Attribute Evauator=default, Search Method=default→Attribute Selection Mode=Use full trining set→ Start
- After getting output for attribute selection save that dataset with new name
- Right click on result list→ Save reduced data

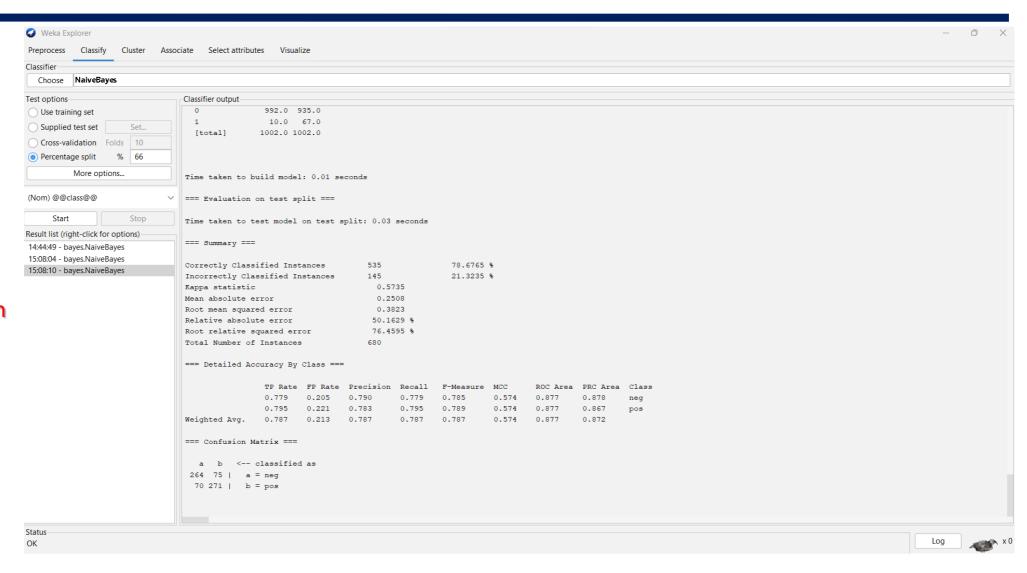Initially we have 1156 attributes, but after attribute selection only 52

# Attribute Selection

# Classification using Naive Bayes Classifier with selected attributes

**Initially**
**Accuracy 77.0588 %**
**for all attributes**

**After attribute Selection**
**Accuracy 78.6765 %**

# Reference

❖ Han, Jiawei Kamber, Micheline Pei and Jian, "Data Mining: Concepts and Techniques",Elsevier Publishers, ISBN:9780123814791, 9780123814807.