

Text Classification and Clustering with WEKA

A guided example by
Sergio Jiménez



The Task

Building a model for movies revisions in English for classifying it into positive or negative.



Sentiment Polarity Dataset Version 2.0

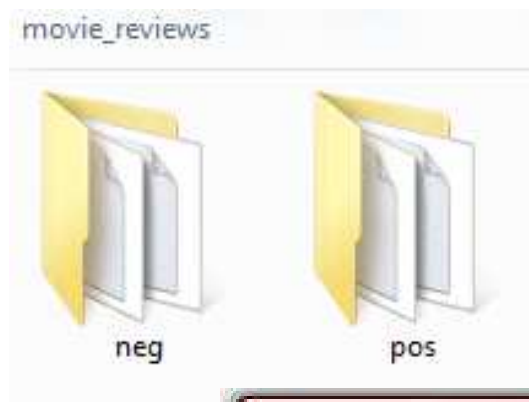
1000 positive movie review and 1000 negative review texts from:

Thumbs up? Sentiment Classification using Machine Learning Techniques. Bo Pang, Lillian Lee, and Shivakumar Vaithyanathan. Proceedings of EMNLP, pp. 79--86, 2002.

“Our data **source** was the **Internet Movie Database** (IMDb) archive of the rec.arts.movies.reviews newsgroup.³ We selected only reviews where the **author rating** was **expressed** either with stars or some **numerical value** (other conventions varied too widely to allow for automatic processing). Ratings were automatically extracted and converted into one of three categories: positive, negative, or neutral. For the work described in this paper, we concentrated **only** on discriminating between **positive** and **negative** sentiment.”

<http://www.cs.cornell.edu/people/pabo/movie-review-data/>

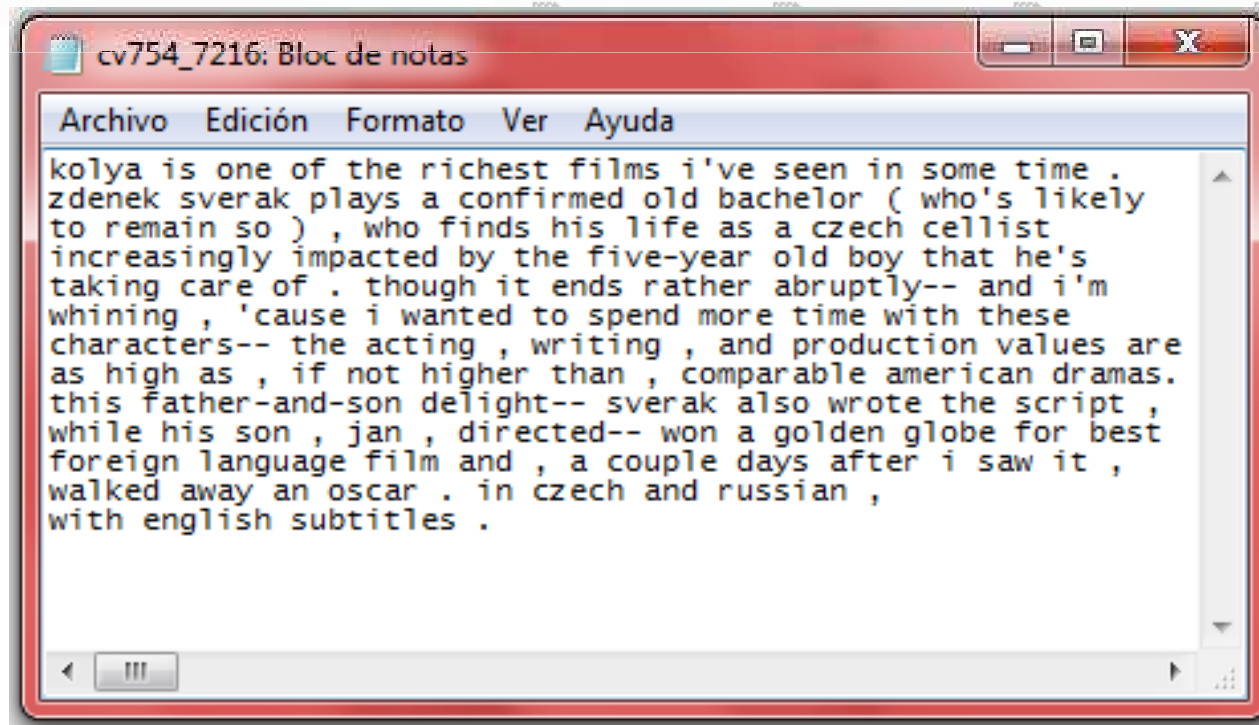
The Data (1/2)



Biblioteca Documentos

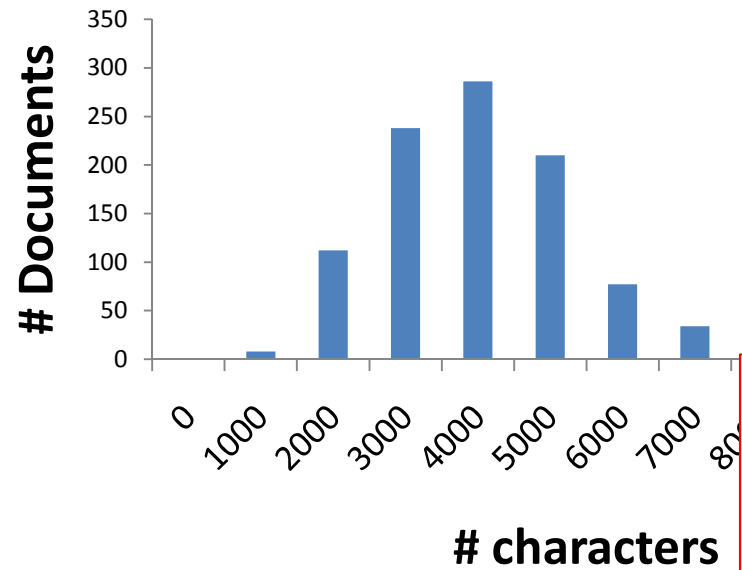
pos

cv754_7216	cv114_18398	cv938_10220	cv013_10159
cv280_8267	cv471_16858	cv424_8831	cv253_10077
cv825_5063	cv230_7428	cv763_14729	cv722_7110
cv057_7453	cv075_6500	cv319_14727	cv082_11080
cv640_5378	cv058_8025	cv430_17351	cv312_29377
			cv361_28944
			cv931_17563
			cv170_3006

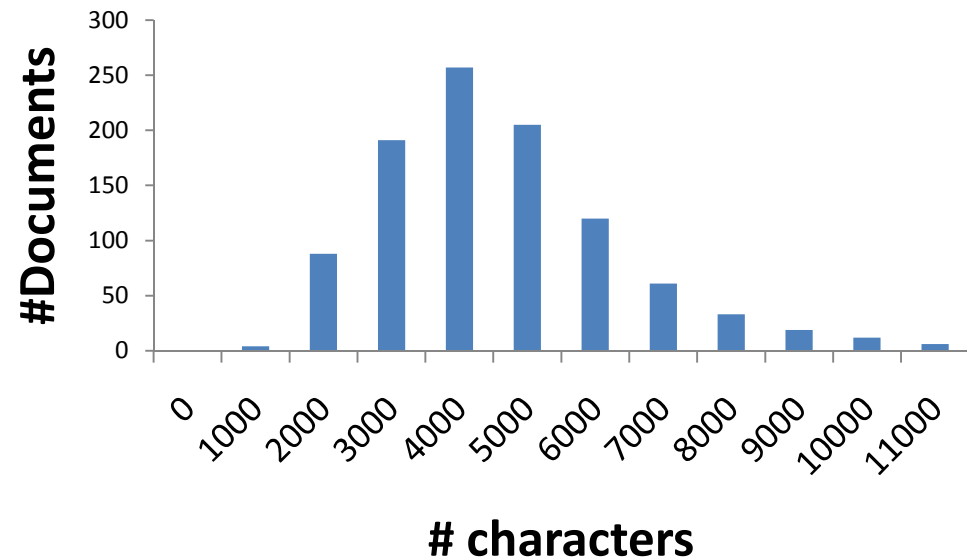


The Data (2/2)

1000 negative revisions histogram



1000 positive revisions histogram



What WEKA is?



- “Weka is a collection of machine learning algorithms for data mining tasks”.
- “Weka contains tools for:
 - data pre-processing,
 - classification,
 - regression,
 - clustering,
 - association rules,
 - and visualization”

Where to start?



WEKA



Rechercher

Environ 969 000 résultats (0,15 secondes)

[Recherche avancée](#)



Tout



Images



Vidéos



Plus

Recherche sur le Web

[Rechercher les pages en français](#)

Date indifférente

2 derniers jours



Plus d'outils

► [Weka 3 - Data Mining with Open Source Machine Learning Software in ...](#) ☆ - 4

visites - 15:23 - [[Traduire cette page](#)]

Collection of machine learning algorithms for solving data mining problems implemented in Java and open sourced under the GPL.

www.cs.waikato.ac.nz/~weka/ - En cache - Pages similaires

[Éditions Weka](#) ☆

Protection sociale des personnels médicaux et hospitaliers www.weka.fr. Maîtrisez les subtilités de chaque situation :... Rémunération et paie des personnels ...

[Marchés Publics](#) - [RH Publiques](#) - [Action sociale](#)

www.weka.fr/ - En cache - Pages similaires

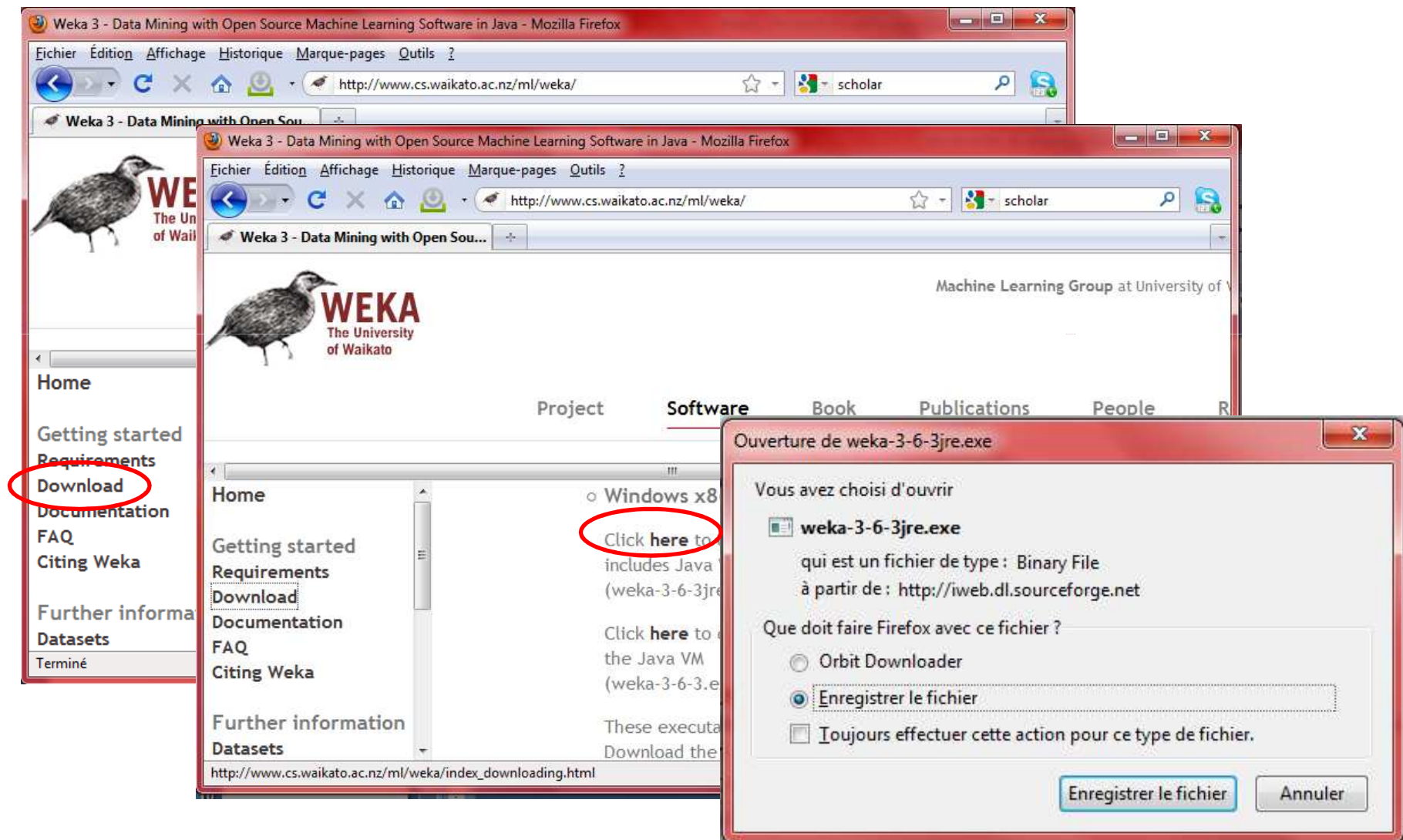
[Weka---Machine Learning Software in Java | Download Weka---Machine ...](#) ☆

- [[Traduire cette page](#)]

Get **Weka**---Machine Learning Software in Java at SourceForge.net. Fast, secure and free downloads from the largest Open Source applications and software ...

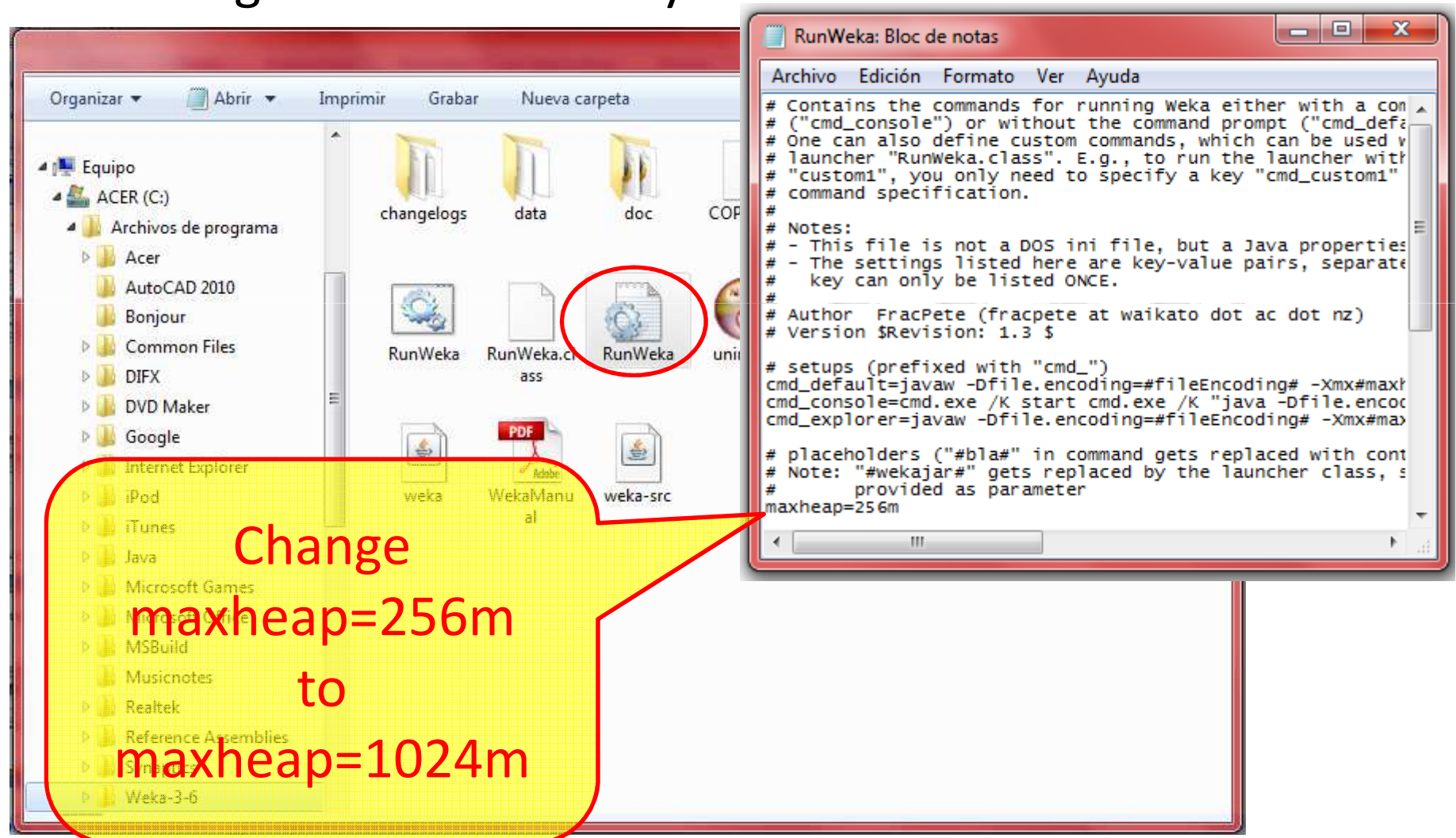
sourceforge.net > [Projects](#) - En cache - Pages similaires

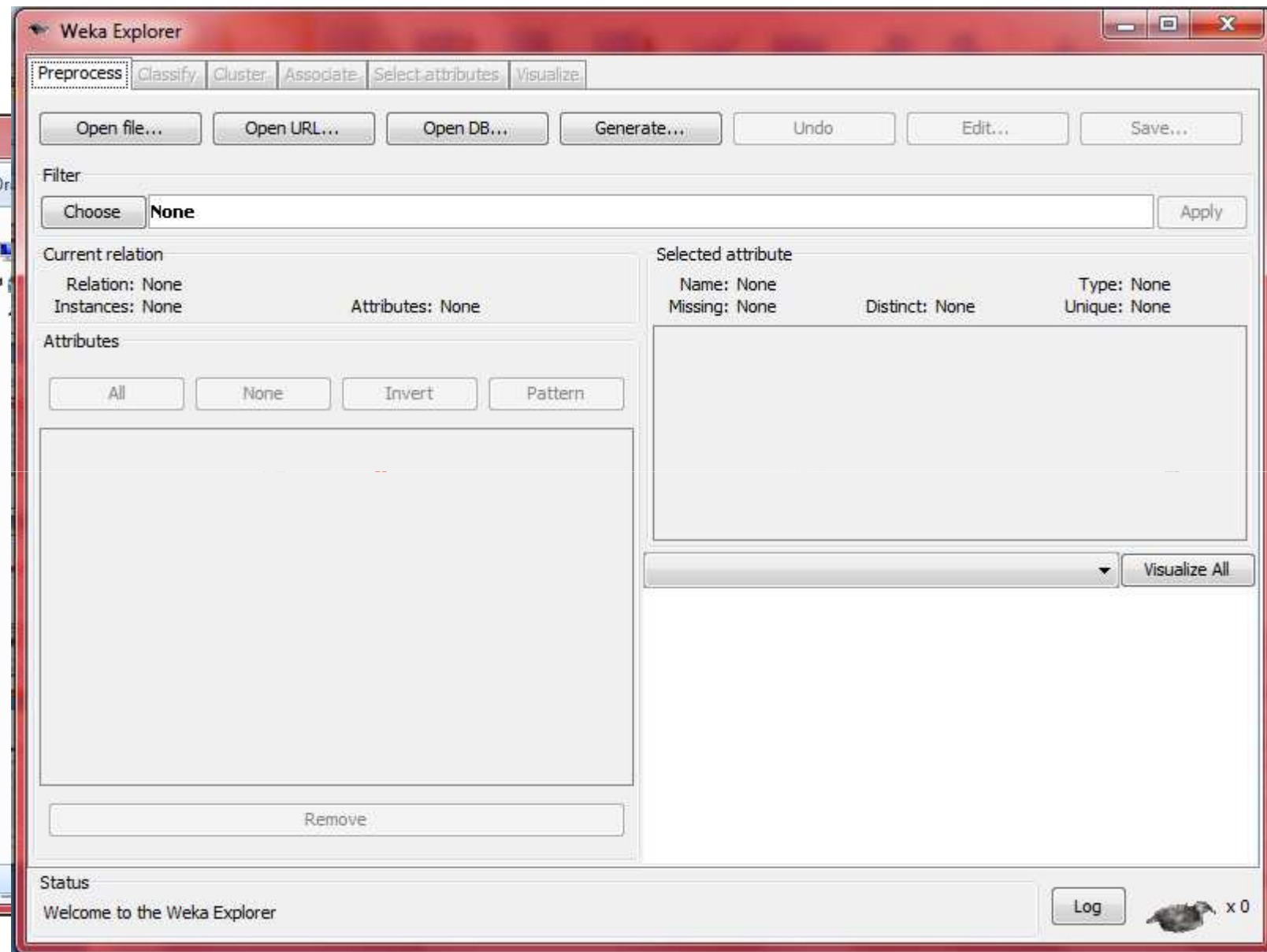
Getting WEKA



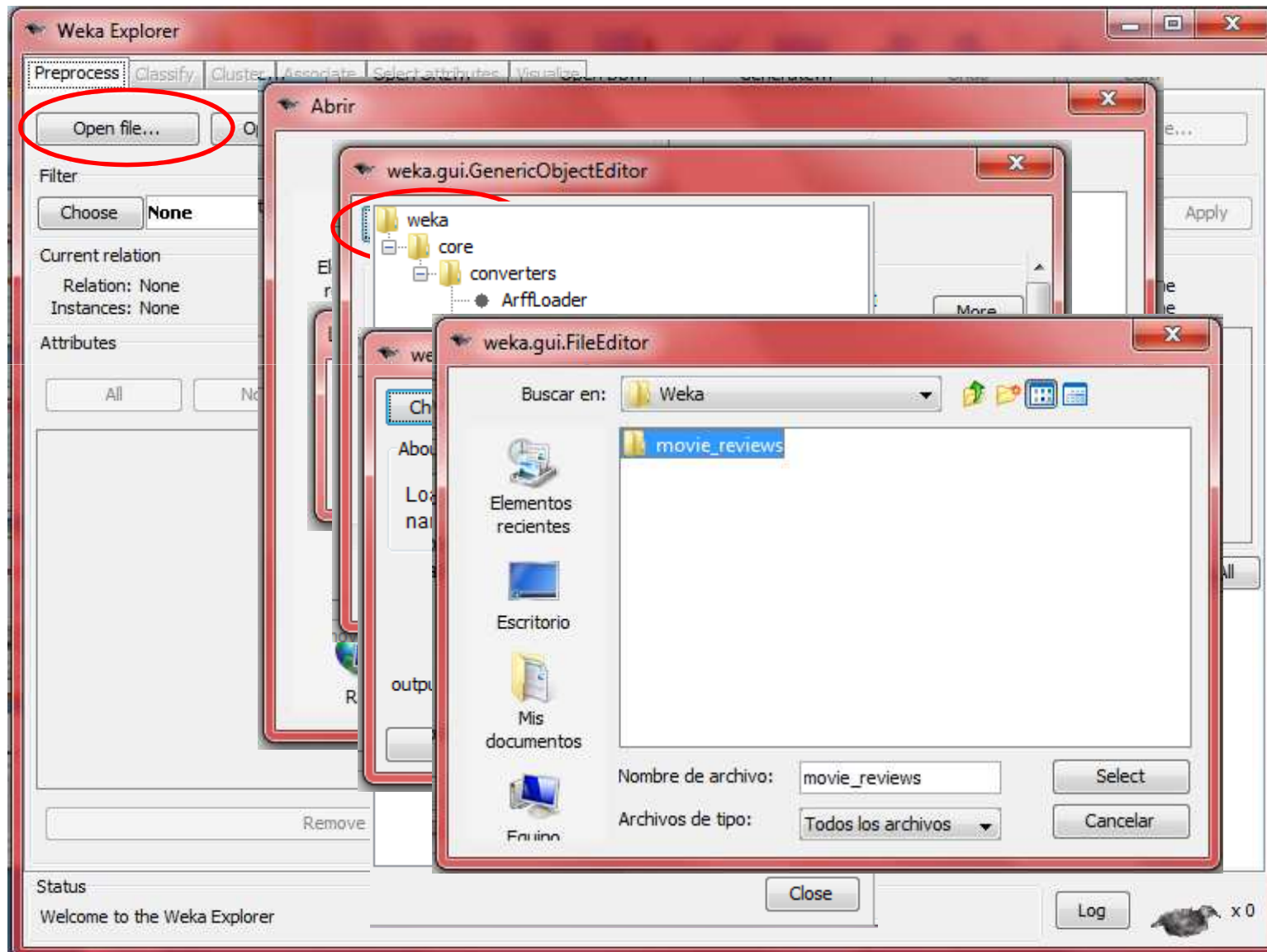
Before Running WEKA

Increasing available memory for Java in RunWeka.ini

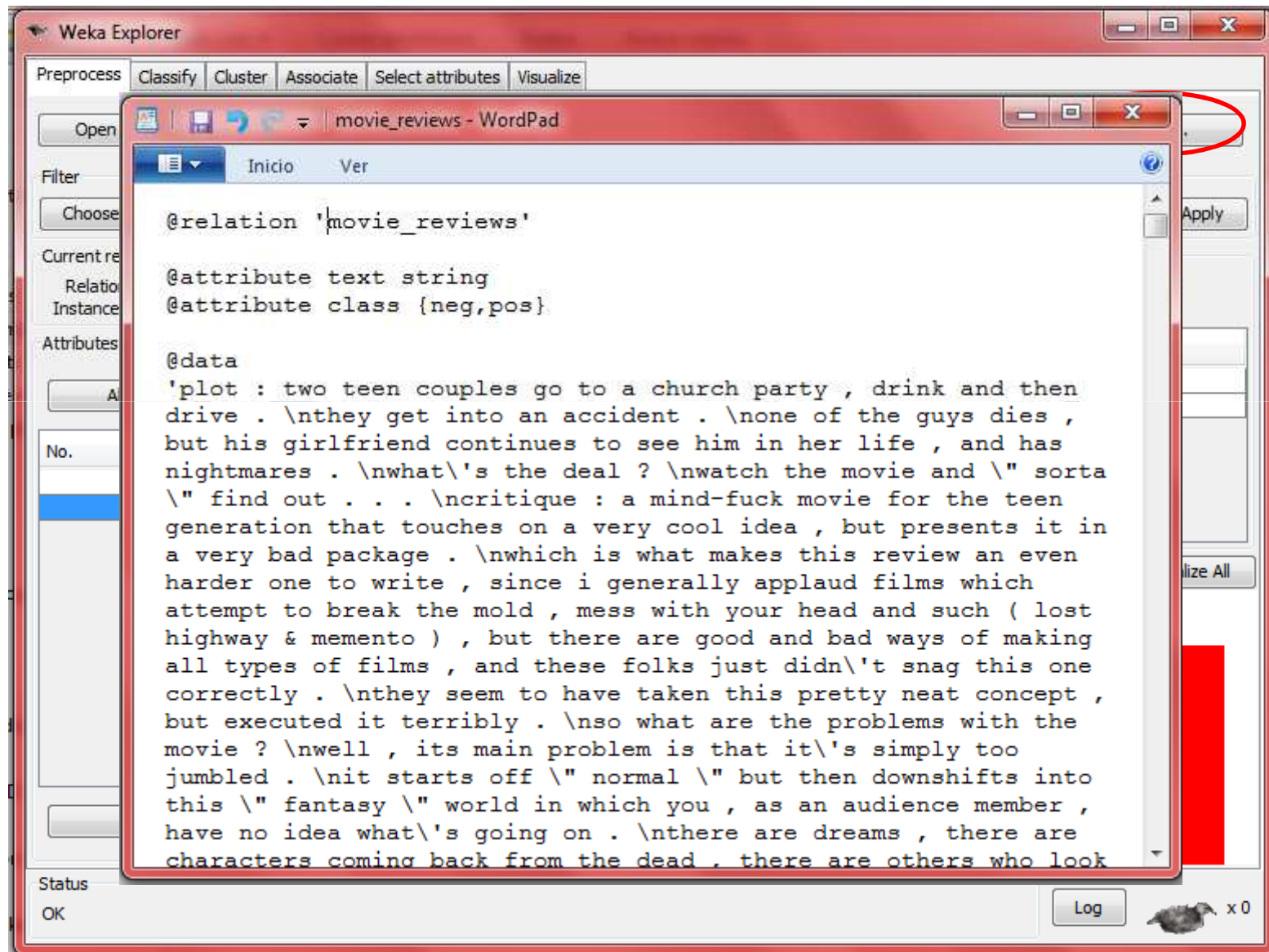




Creating a .arff dataset



Saving the .arff dataset



From text to vectors

$$V = [v_1, v_2, v_3, \dots, v_n, class]$$

review₁ = “great movie” review₃ = “worst film ever”

review₂ = “excellent film” review₄ = “sucks”

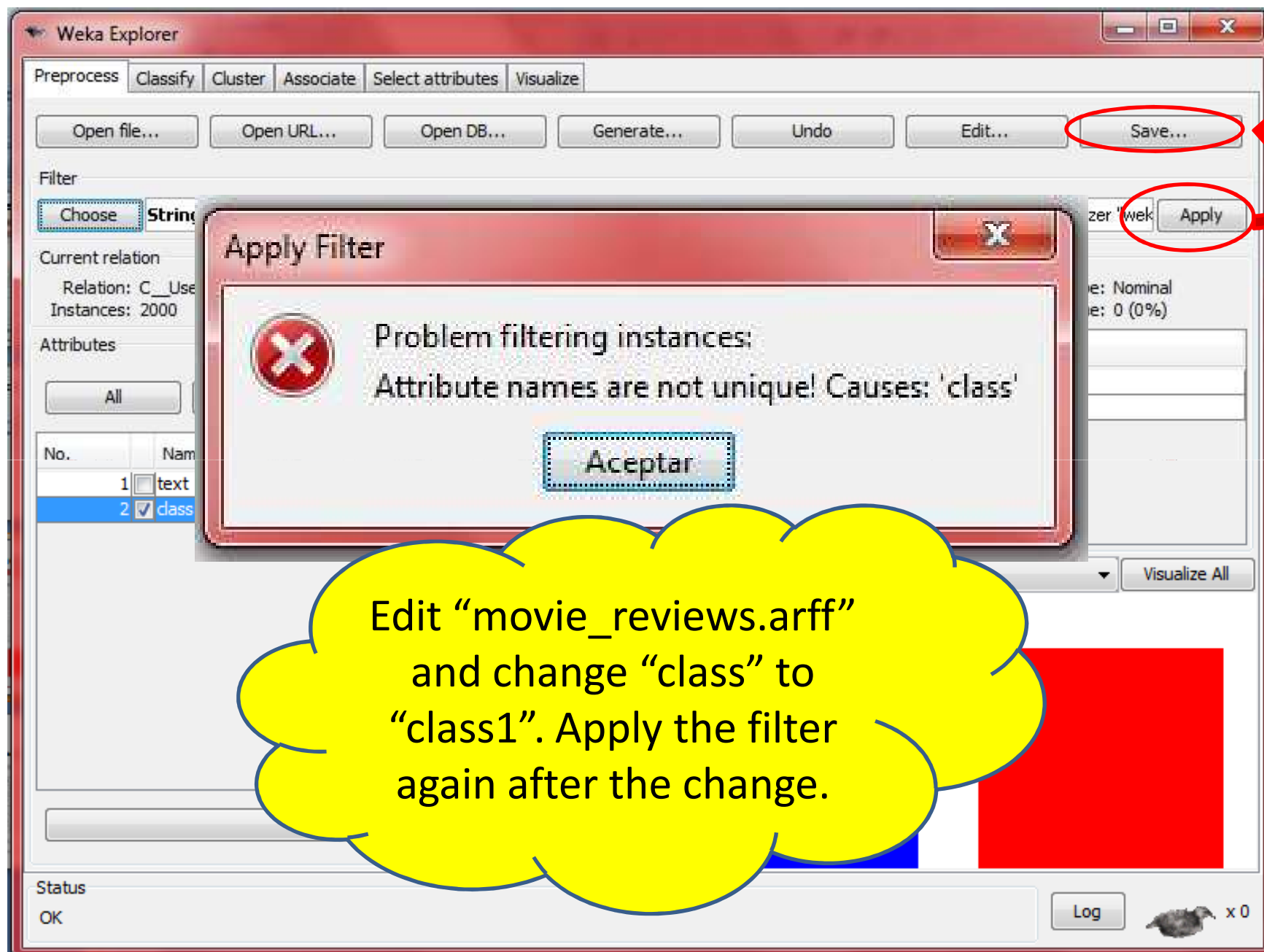
ever
excellent
film
great
movie
sucks
worst

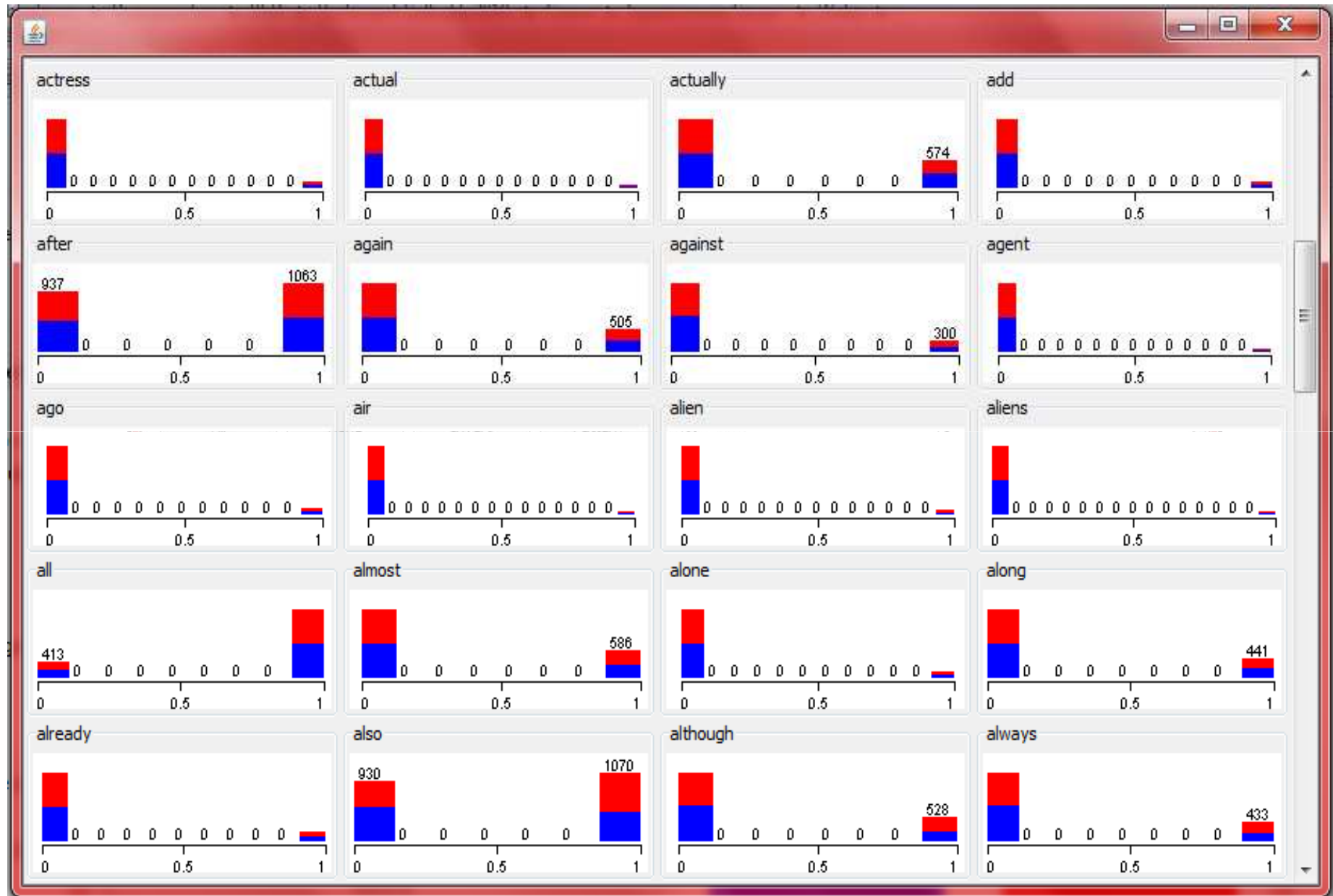
$$V_1 = [0, 0, 0, 1, 1, 0, 0, +]$$

$$V_2 = [0, 1, 1, 0, 0, 0, 0, +]$$

$$V_3 = [1, 0, 1, 0, 0, 0, 1, -]$$

$$V_4 = [0, 0, 0, 0, 0, 1, 0, -]$$





StringToWordVector

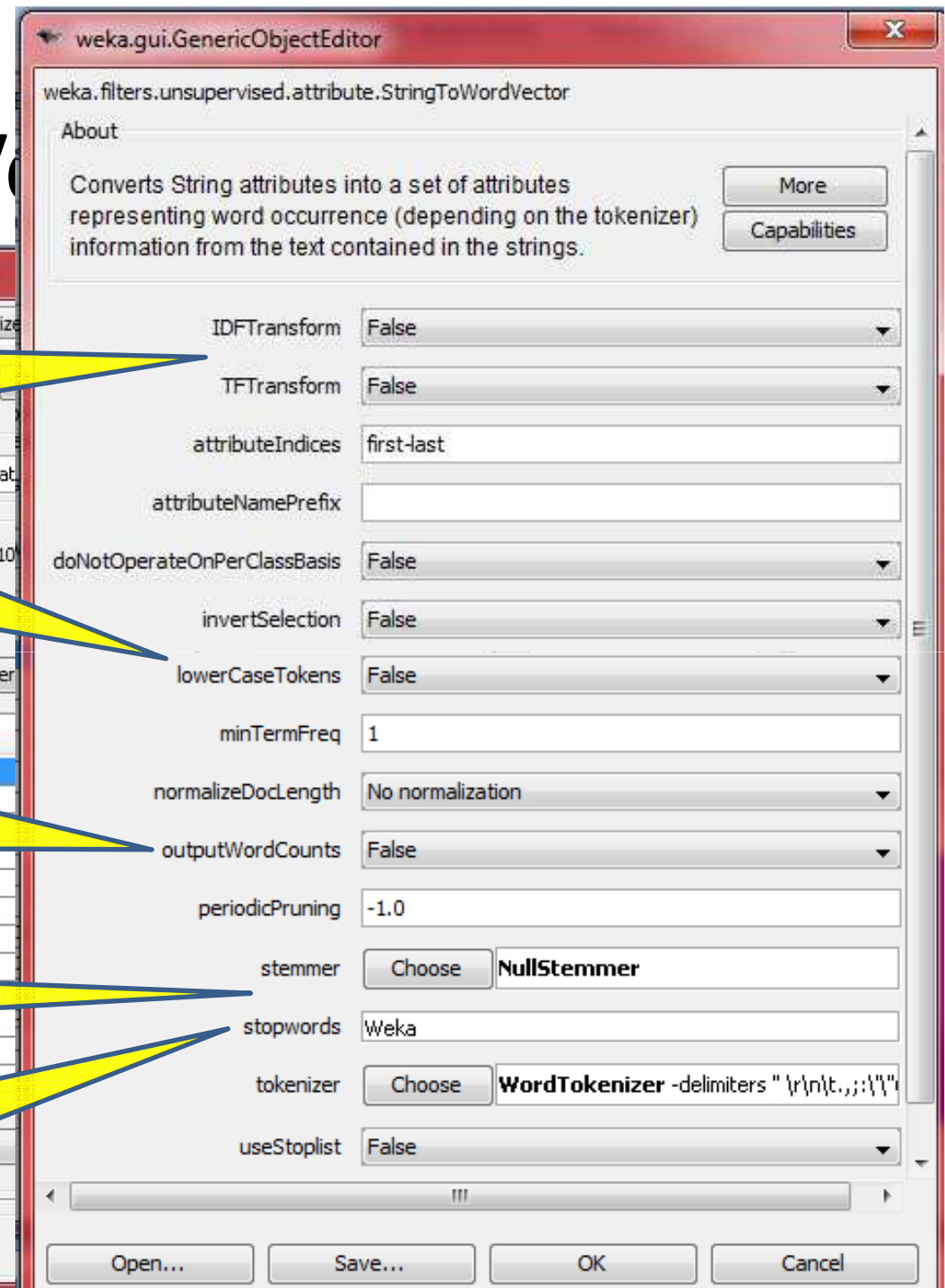
TF-IDF weighing

lowerCase conversion

Use frequencies instead of
single presence

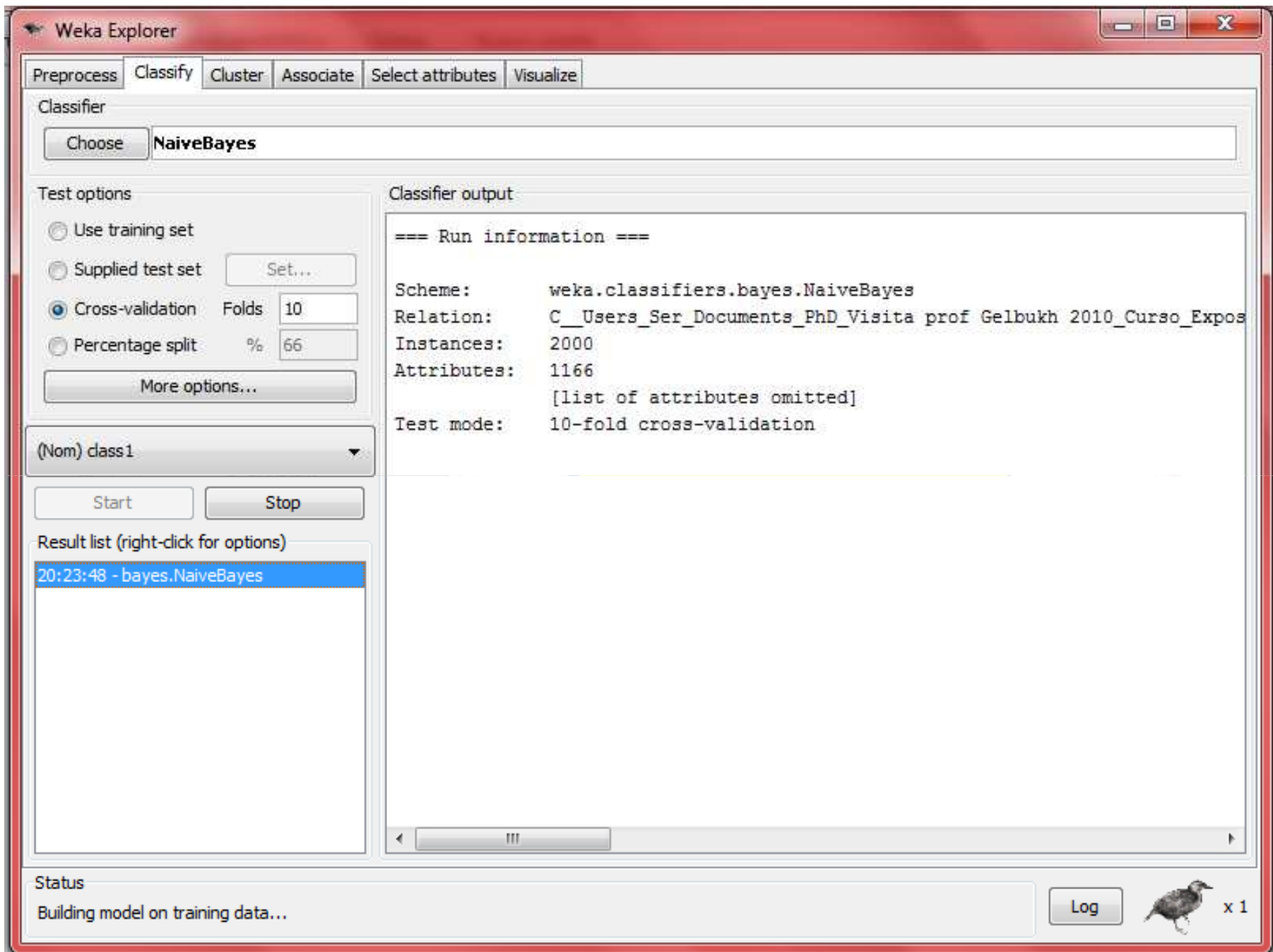
Stemming

Stopwords removal using a list
of words in a file



Generating datasets for experiments

<i>dataset file name</i>	<i>Stopwords</i>	<i>Stemming</i>	<i>Presence or freq.</i>
movie_reviews_1.arff		no	presence
movie_reviews_2.arff		no	frequency
movie_reviews_3.arff		yes	presence
movie_reviews_4.arff		yes	frequency
movie_reviews_5.arff	removed	no	presence
movie_reviews_6.arff	removed	no	frequency
movie_reviews_7.arff	removed	yes	presence
movie_reviews_8.arff	removed	yes	frequency



Results

Correctly Classified Instances	1616	80.8	%
Incorrectly Classified Instances	384	19.2	%
Kappa statistic	0.616		
Mean absolute error	0.1918		
Root mean squared error	0.4111		
Relative absolute error	38.3507	%	
Root relative squared error	82.2217	%	
Total Number of Instances	2000		

=== Detailed Accuracy By Class ===

	TP Rate	FP Rate	Precision	Recall	F-Measure	ROC Area
	0.832	0.216	0.794	0.832	0.813	0.897
	0.784	0.168	0.824	0.784	0.803	0.897
Weighted Avg.	0.808	0.192	0.809	0.808	0.808	0.897

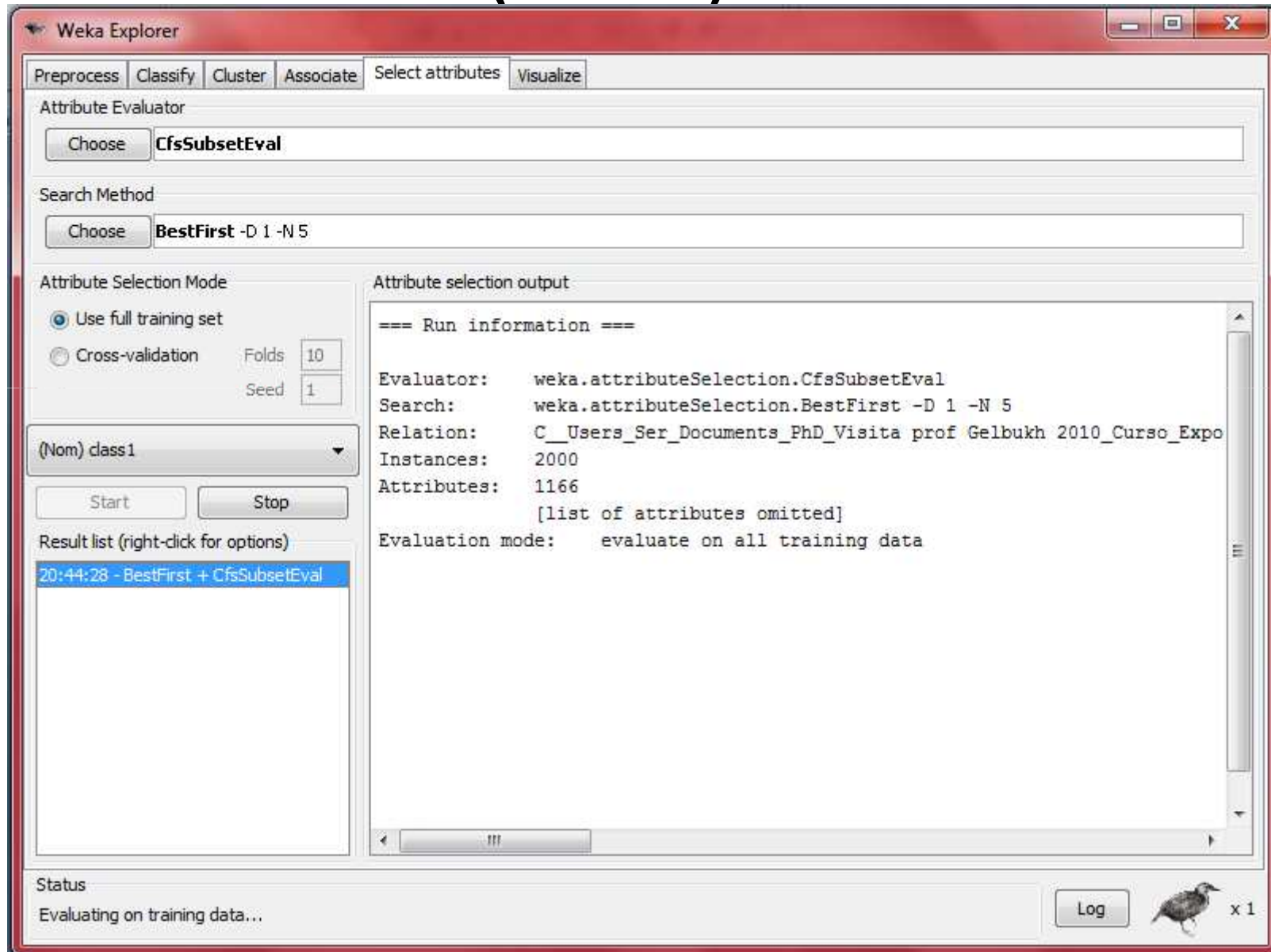
=== Confusion Matrix ===

a	b	<-- classified as
832	168	a = neg
216	784	b = pos

Results Correctly Classified Reviews

<i>dataset name</i>	<i>Stopwords</i>	<i>Stemming</i>	<i>Presence or freq.</i>	<i>Naive Bayes 3- fold</i>	<i>NaiveBayes Multinomial 3-fold</i>
movie_reviews_1.arff		no	presence	80.65%	83.80%
movie_reviews_2.arff		no	frequency	69.30%	78.65%
movie_reviews_3.arff		yes	presence	79.40%	82.15%
movie_reviews_4.arff		yes	frequency	68.10%	79.70%
movie_reviews_5.arff	removed	no	presence	81.80%	84.35%
movie_reviews_6.arff	removed	no	frequency	69.40%	81.75%
movie_reviews_7.arff	removed	yes	presence	78.90%	82.40%
movie_reviews_8.arff	removed	yes	frequency	68.30%	80.50%

Attribute (word) Selection



Selected Attributes (words)

also
awful
bad
boring
both
dull
fails
great
joke
lame
life
many
maybe
mess
nothing
others
perfect
performances

pointless
poor
ridiculous
script
seagal
sometimes
stupid
tale
terrible
true
visual
waste
wasted
world
worst
animation
definitely

deserves



wonderfully

Pruned movie_reviews_1.arff dataset

Weka Explorer

Preprocess | Classify | Cluster | Associate | Select attributes | Visualize

Open file... Open URL... Open DB... Generate... Undo Edit... Save...

Filter
Choose **None** Apply

Current relation
Relation: C:\Users\Ser_Documents_PhD_Visita_prof_Gelbukh_2010_C...
Instances: 2000 Attributes: 51

Attributes
All None Invert Pattern

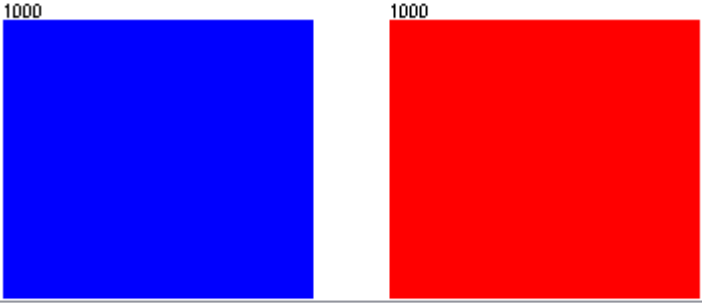
No.	Name
1	<input checked="" type="checkbox"/> class1
2	<input type="checkbox"/> also
3	<input type="checkbox"/> awful
4	<input type="checkbox"/> bad
5	<input type="checkbox"/> boring
6	<input type="checkbox"/> both
7	<input type="checkbox"/> dull
8	<input type="checkbox"/> fails
9	<input type="checkbox"/> great
10	<input type="checkbox"/> joke
11	<input type="checkbox"/> lame
12	<input type="checkbox"/> life
13	<input type="checkbox"/> many

Remove

Selected attribute
Name: class1
Missing: 0 (0%)
Distinct: 2
Type: Nominal
Unique: 0 (0%)

No.	Label	Count
1	neg	1000
2	pos	1000

Class: class1 (Nom) Visualize All



Status
OK Log x 0

Naïve Bayes with the pruned dataset

The screenshot shows the Weka Explorer application window. The 'Classify' tab is active, and the 'NaiveBayes' classifier is selected. The 'Test options' section shows 'Cross-validation' with 3 folds. The 'Result list' on the left shows the current run at 21:10:00. The 'Classifier output' pane displays the following summary statistics:

Metric	Value	Percentage
Correctly Classified Instances	1621	81.05 %
Incorrectly Classified Instances	379	18.95 %
Kappa statistic	0.621	
Mean absolute error	0.2157	
Root mean squared error	0.3814	81.05 %
Relative absolute error	43.1437 %	
Root relative squared error	76.276 %	18.95 %
Total Number of Instances	2000	

Below the summary, the 'Detailed Accuracy By Class' table is shown:

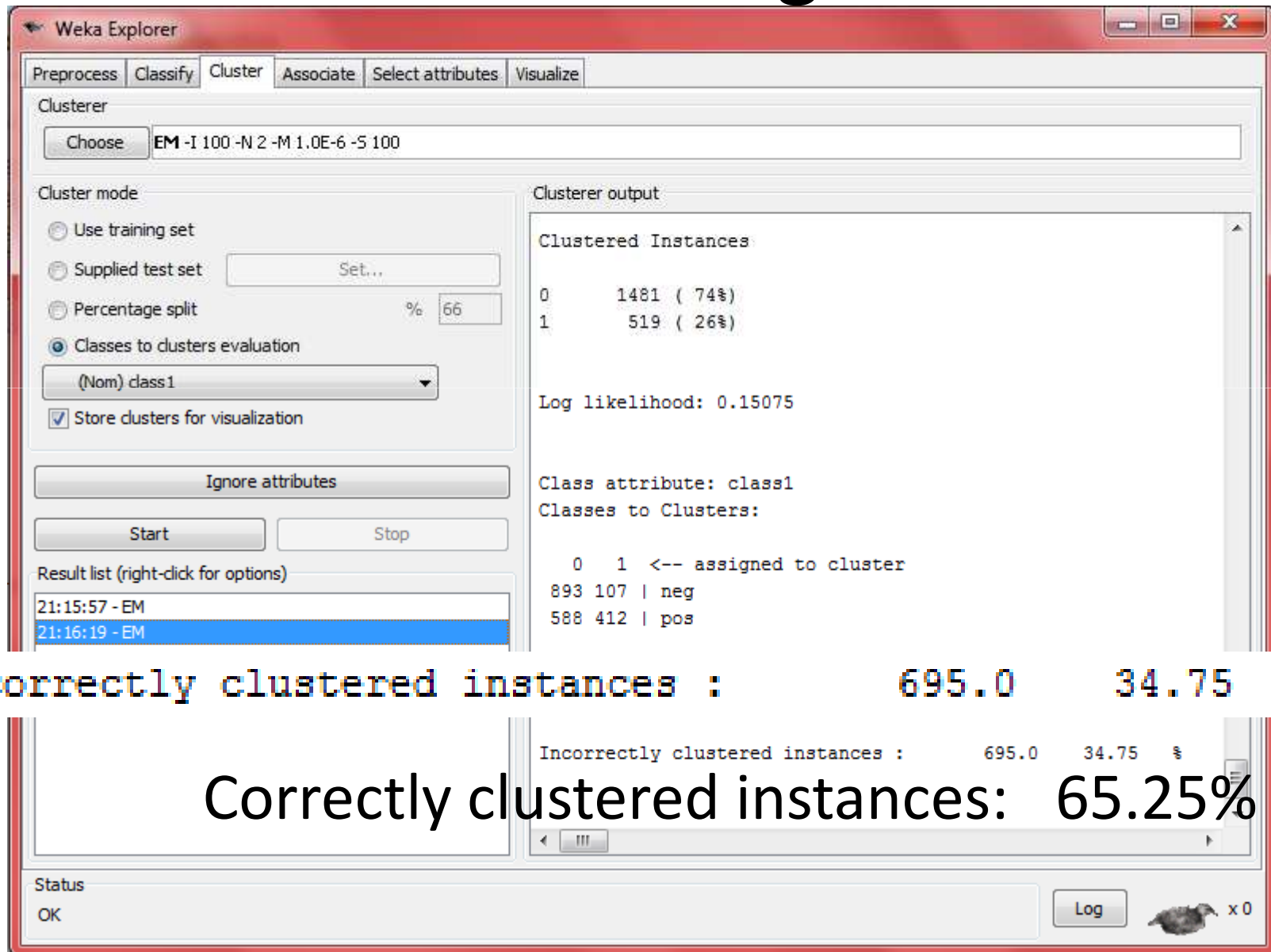
	TP Rate	FP Rate	Precision	Recall	F-Measure	ROC Area
Weighted Avg.	0.811	0.19	0.811	0.811	0.81	0.888

The 'Confusion Matrix' section shows the following data:

```
a  b  <-- classified as
804 196 | a = neg
183 817 | b = pos
```

The status bar at the bottom indicates 'Status OK' and includes a 'Log' button.

Clustering



The screenshot shows the Weka Explorer interface with the 'Cluster' tab selected. The 'Clusterer' dropdown is set to 'EM -I 100 -N 2 -M 1.0E-6 -S 100'. The 'Cluster mode' section has 'Use training set' selected, and 'Classes to clusters evaluation' is checked with '(Nom) class1' selected. The 'Clusterer output' pane displays the following results:

```
Clustered Instances
0      1481 ( 74%)
1       519 ( 26%)

Log likelihood: 0.15075

Class attribute: class1
Classes to Clusters:

  0   1  <-- assigned to cluster
893 107 | neg
588 412 | pos
```

Below the output, a text overlay reads: 'Incorrectly clustered instances : 695.0 34.75 %'. The 'Result list' at the bottom shows two entries: '21:15:57 - EM' and '21:16:19 - EM', with the latter selected. The status bar at the bottom indicates 'OK' and 'Log'.

Other results

Results of Pang et al. (2002) with version 1.0 of the dataset with 700+ and 700-

	Features	# of features	frequency or presence?	NB	ME	SVM
(1)	unigrams	16165	freq.	78.7	N/A	72.8
(2)	unigrams	”	pres.	81.0	80.4	82.9
(3)	unigrams+bigrams	32330	pres.	80.6	80.8	82.7
(4)	bigrams	16165	pres.	77.3	77.4	77.1
(5)	unigrams+POS	16695	pres.	81.5	80.4	81.9
(6)	adjectives	2633	pres.	77.0	77.7	75.1
(7)	top 2633 unigrams	2633	pres.	80.3	81.0	81.4
(8)	unigrams+position	22430	pres.	81.0	80.1	81.6

No stemming or stoplists were used.

are the average three-fold cross-validation results

Thanks