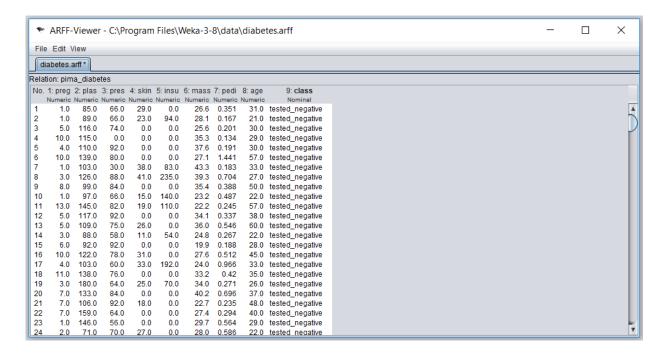# Classification using Weka & Scikit-learn

*PROF. K. H. LIM*

50.038 Computational Data Science

# Load in datasets

o Weka already has some pre-installed datasets
  ◦ Go to Tools → ArffViewer → File → Open
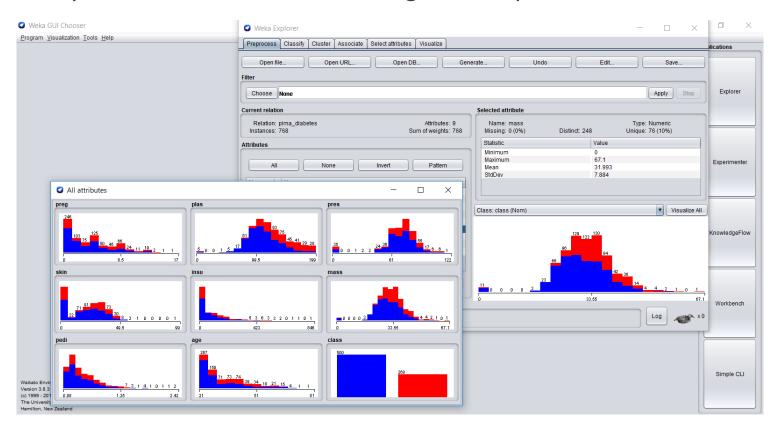  ◦ Then select dataset ./Weka-3-8/data/diabetes.arff

# Attributes in diabetes.arff

o **Preg**: Number of times pregnant

o **Plas**: Plasma glucose concentration

o **Pres**: Diastolic blood pressure (mm Hg)

o **Skin**: Triceps skin fold thickness (mm)

o **Insu**: 2-Hour serum insulin (mu U/ml)

o **Mass**: Body mass index (weight in kg/(height in m)^2)

o **Pedi**: Diabetes pedigree function

o **Age**: Age (years)
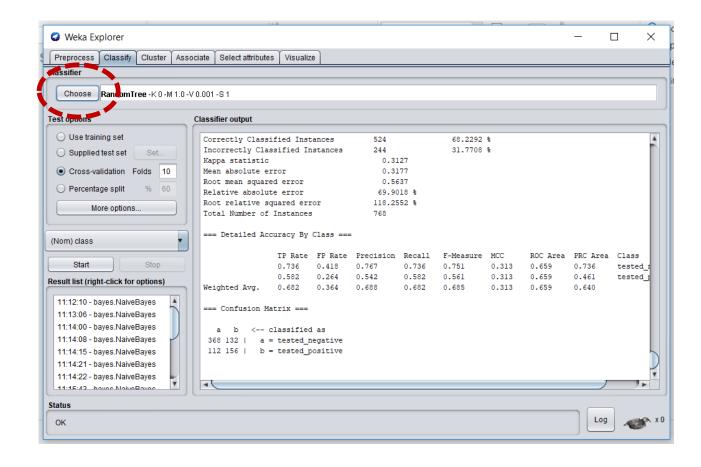
o **Class**: Test results for diabetes

# Weka Data Explorer

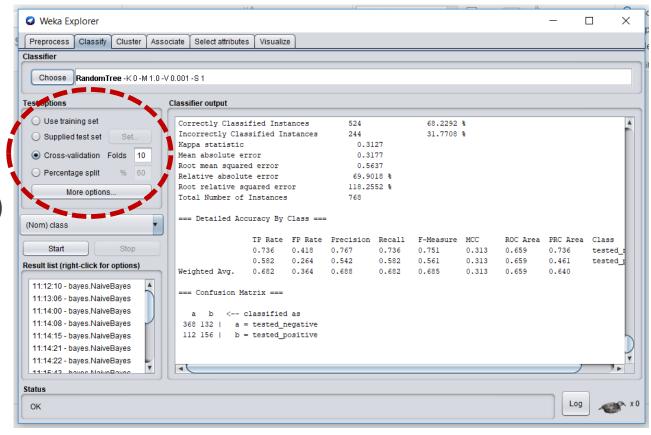o Study the diabetes.arff dataset using Weka Explorer

# Training and Testing Classifiers

o Select from multiple classifiers (e.g., Naïve Bayes, Decision Trees, Random Forests, etc)

# Training and Testing Classifiers

○ Various evaluation approach (e.g., using pre-fixed train/test sets, k-fold cross validation, etc)

# Exercise

1. Load in the diabetes.arff dataset

2. How many features does this dataset have? What will you use as labels? 8

3. Use Weka explorer, find out: (i) what is the average "mass"; and (ii) what are the minimum and maximum "age" 31.993
   21    81

4. Run the "RandomTree" and "RandomForest" classifiers using a 10-fold cross validation.
   ◦ What do you observe in terms of their performance and running time?

# Python Scikit-learn

o Machine learning library based on Python

o Contains functionalities for data pre-processing, classification, clustering, etc

# Twitter Dataset

o Using a Twitter sentiment dataset
   ◦ http://cs.stanford.edu/people/alecmgo/trainingandtestdata.zip
   ◦ https://docs.google.com/file/d/0B04GJPshIjmPRnZManQwWEdTZjg/edit

o This dataset comprises 1.6M tweets with various columns, we will make use of the first and last column (sentiment label and tweet text)
   ◦ For sentiment, a value of 4 = positive and 0 = negative

# Load Packages

o Import relevant packages

```python
1   # load in required packages
2   from sklearn.feature_extraction.text import CountVectorizer
3   from sklearn.naive_bayes import MultinomialNB
4   from sklearn.pipeline import Pipeline
5   from sklearn.model_selection import train_test_split
6   from sklearn.preprocessing import FunctionTransformer
7   from sklearn import metrics
8   from sklearn.metrics import accuracy_score
9   import pandas as pd
10  import numpy as np
```

# Load Dataset

o Does dataset contain headers?

o What type of encoding for text?

```
1  # our dataset has no column names, so we define them
2  colnames = ['label','id','date','query','user','text']
3
4  # load in training/test set
5  df_train = pd.read_csv('training.1600000.processed.noemoticon.csv',
6                         header=None, names=colnames, encoding='windows-1252')
7  df_test = pd.read_csv('testdata.manual.2009.06.14.csv',
8                         header=None, names=colnames, encoding='windows-1252')
```

# Sanity Check

○ Is the dataset how you expected it to be?

```
In [3]:    1  # check dimensions
           2  df_train.shape
```

```
Out[3]: (1600000, 6)
```

```
In [4]:    1  # check that training set is correct
           2  df_train.head()
```

Out[4]:

|   | label | id | date | query | user | text |
|---|-------|----|------|-------|------|------|
| 0 | 0 | 1467810369 | Mon Apr 06 22:19:45 PDT 2009 | NO_QUERY | _TheSpecialOne_ | @switchfoot http://twitpic.com/2y1zl - Awww, t... |
| 1 | 0 | 1467810672 | Mon Apr 06 22:19:49 PDT 2009 | NO_QUERY | scotthamilton | is upset that he can't update his Facebook by ... |
| 2 | 0 | 1467810917 | Mon Apr 06 22:19:53 PDT 2009 | NO_QUERY | mattycus | @Kenichan I dived many times for the ball. Man... |
| 3 | 0 | 1467811184 | Mon Apr 06 22:19:57 PDT 2009 | NO_QUERY | ElleCTF | my whole body feels itchy and like its on fire |
| 4 | 0 | 1467811193 | Mon Apr 06 22:19:57 PDT 2009 | NO_QUERY | Karoli | @nationwideclass no, it's not behaving at all.... |

```
In [5]:    1  # check the class distribution
           2  df_train['label'].value_counts()
```

```
Out[5]: 4    800000
        0    800000
        Name: label, dtype: int64
```

# Bag of Words (BoW)

o Need to count the frequency of each word

o Use the CountVectorizer() function

```
In [10]:    1  # counting of words (bag of words)
            2  bowVect = CountVectorizer()
            3  trainBow = bowVect.fit_transform(df_train.text)
            4  trainBow.shape

Out[10]:    (1600000, 685256)
```

# Training and Testing

o Convert test set to BoW representation

o Train our model by calling fit() on our training set and labels

o Test our model by calling predict()

```
1  # train and test our model
2  testBow = bowVect.transform(df_test.text) # not fit_transform()
3  mnbClf = MultinomialNB().fit(trainBow, train_labels)
4  predicted = mnbClf.predict(testBow)
```

# Evaluation Scores

o Look at Accuracy, Precision, Recall and F-measure

```python
# evaluating our model
print(metrics.classification_report(test_labels, predicted))
print(accuracy_score(test_labels, predicted))
```

In [12]:

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.79 | 0.82 | 0.81 | 177 |
| 4 | 0.82 | 0.79 | 0.81 | 182 |
| avg / total | 0.81 | 0.81 | 0.81 | 359 |

0.807799442896936

# Different Features and Models

o What happens if we need to evaluate different features or classifiers?

◦ Instead of words (uni-grams), use bi-grams and tri-grams?

◦ Instead of Naïve Bayes, use Decision Tree, Random Forest, SVM, etc?


o Can use Scikit-learn's Pipeline

◦ It takes an input, does a series of pre-processing step before feeding it to a final classifier

◦ More information at http://scikit-learn.org/stable/modules/generated/sklearn.pipeline.Pipeline.html

# Exercise

5. Load in Twitter dataset (both files) and count the frequency of each class

6. Implement a simple pipeline with Naïve Bayes classifier with words (unigrams) as features
   ◦ Using "training.1600000.processed.noemoticon.csv", evaluate your model with 80% training and 20% test
   ◦ hint: look at the train_test_split() function

7. Enhance the Naïve Bayes model by removing stop words, and experimenting with bi-grams and tri-grams
   ◦ hint: look at the various parameters for CountVectorizer()