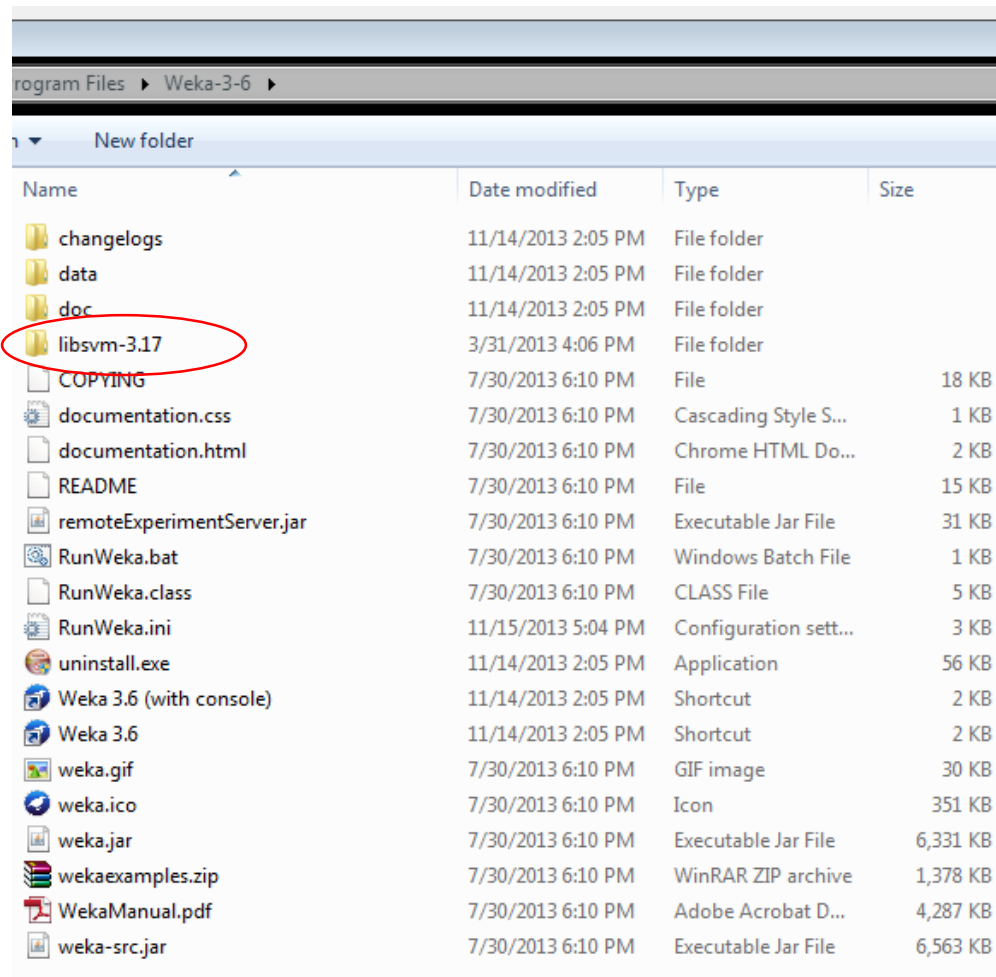For homework 2 and 3, you will be running different machine learning algorithms against the datasets that will be provided to you.

Tools

**WEKA** is a useful machine learning work bench that allows you to quickly run various machine learning algorithms against your datasets. We use this because you do not need to implement (program) the various algorithms for your homework assignments. For this class, we will use the latest version of WEKA available here: http://www.cs.waikato.ac.nz/ml/weka/downloading.html

**LightSIDE** is another machine learning program and we will use it to generate ARFF files which WEKA will read. This is the place where you will be able to determine which features you are going to pass into WEKA for the model building process. We are going to be using an older version of lightSIDE for the purposes of these 2 homework assignments for consistency. <lightsideLink>

**libSVM** is a support vector machine library add-on for Weka. We will be using it for this class. We will bundle it with DNDW but you have to **install it manually by copying it into your Program Files/Weka-3-6/** folder.



Next, edit your RunWeka.ini file and make the following modification:

*cp=%CLASSPATH%;C:/Program Files/Weka-3-6/libsvm-3.17/java/libsvm.jar*

This will add the libsvm libraries to the Classpath used by Weka.

You may also take this time to change this parameter:

*maxheap=2048M*

to something like m*axheap=4096M*

This doubles the memory that Weka can use and will speed things up. Only do this if you have more than 4GB of RAM to devote to running Weka.

**DNDW** (Drag & Drop Weka) is a set of windows batch scripts that automate the laborious process of running Weka experiments and collecting experiment results. These were written by a former student of the class (now current TA) to help save time on homework assignments 2 & 3. Due to the lack of time resources and hardware (TA doesn't have a mac), we are unable to offer a native mac version of the scripts. If you have a mac and would like to port the scripts to run effectively on a mac, please email the TA. Otherwise, this document will offer suggestions on what to do if you have a mac.

Goals

The goal of the homework is to get you thinking about feature selection and identify generally what features work well with certain kinds of algorithms. We also want you to be exposed to different datasets and demonstrate that there is no "winning" algorithm that works on all datasets.
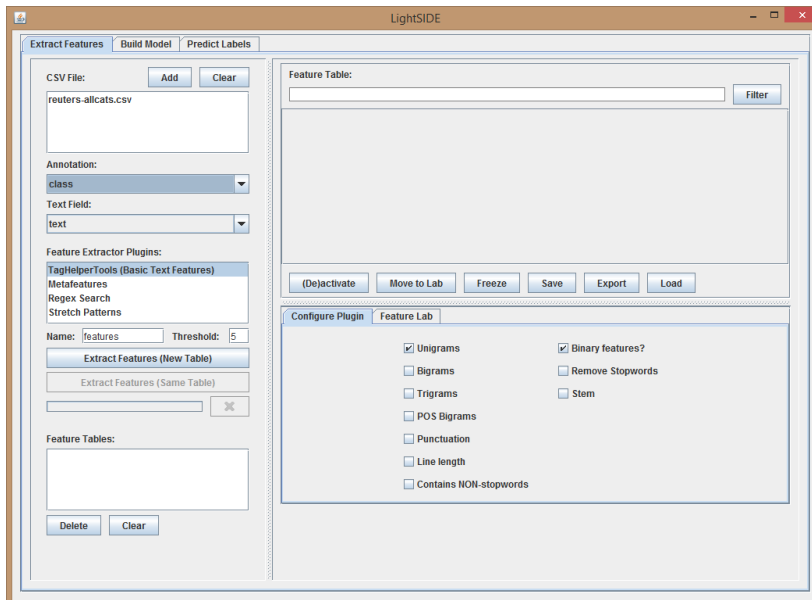
Real world tasks

Keep in mind that there was a lot of pre-processing done on the datasets you are receiving so that they are in formats that can be easily accessed by LightSIDE and WEKA. A real-world problem would involve much more of that labor intensive work and you would probably have to write your own implementations of the various machine learning algorithms as well. However, in this class we have relieved you of those duties ☺ so that you can focus on learning the concepts behind Text Analytics.
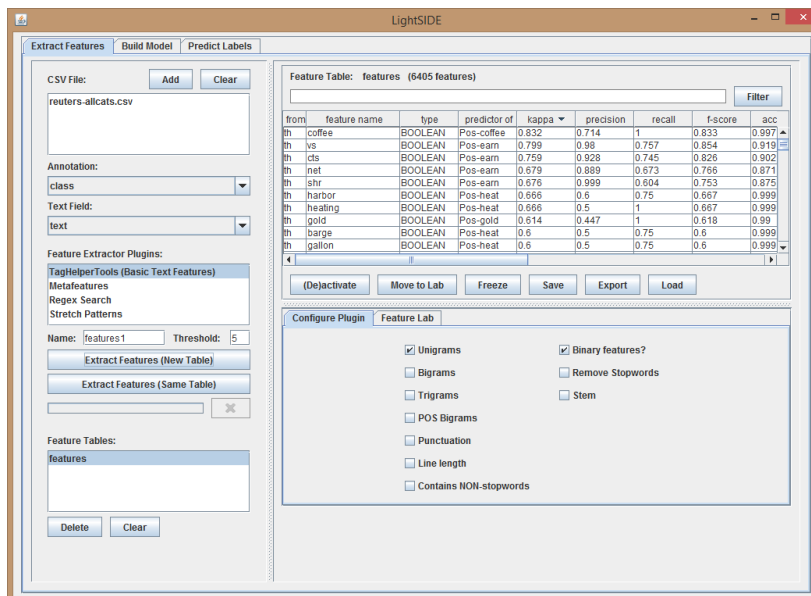
**PC & Mac Users**

1. Download & install WEKA (PC & Mac)
2. Download & install older version of LightSide (extract it anywhere you want) (PC & Mac)
3. Download & extract DNDW (put DNDW anywhere you want) (PC Only)
4. Copy libSVM folder into Program Files/Weka-3-6/ (PC Only) This means the path Program Files/Weka-3-6/libsvm-3.17/java/libsvm.jar is valid.

1 should be rather trivial. For 2 & 3, you can put it anywhere as the batch scripts will handle the file paths.
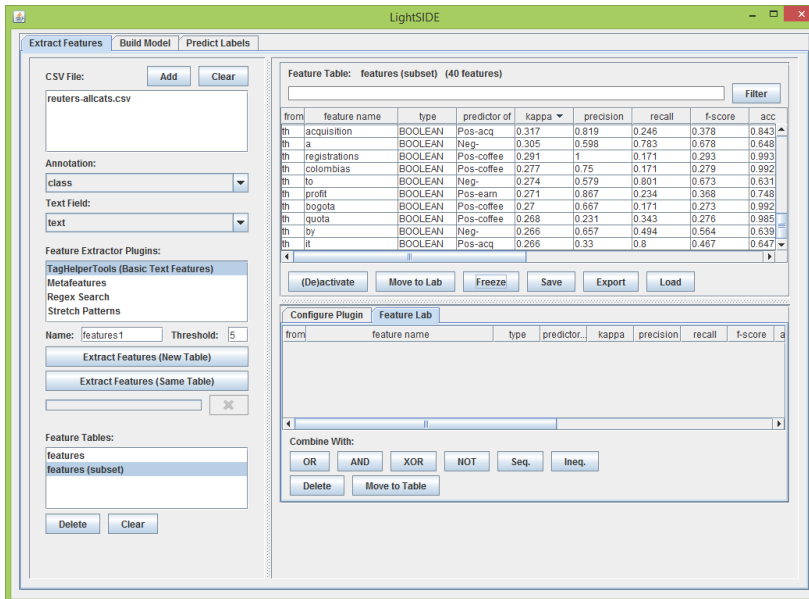
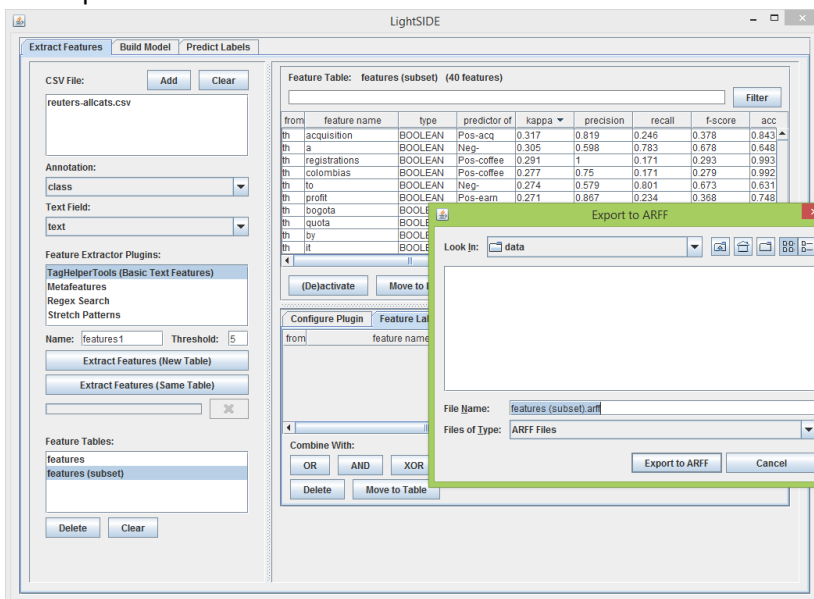5. Run LightSide and add a CSV file (reuters-allcats.csv in this example)



6. A bunch of features pop up according to the threshold and features (Unigrams, Binary) we selected.

7. For the homework a typical task is to get top 40 counts by kappa. The lightSide interface is a little clumsy for doing that. Notice that LightSide sorts kappa for you in descending order. In order to get top 40, we have to skip the first 40 features, select the 41$^{st}$ features, scroll till the last feature, HOLD DOWN SHIFT KEY & click on the last feature – highlighting 41$^{st}$ to last features, click on deactivate, and then click on freeze.

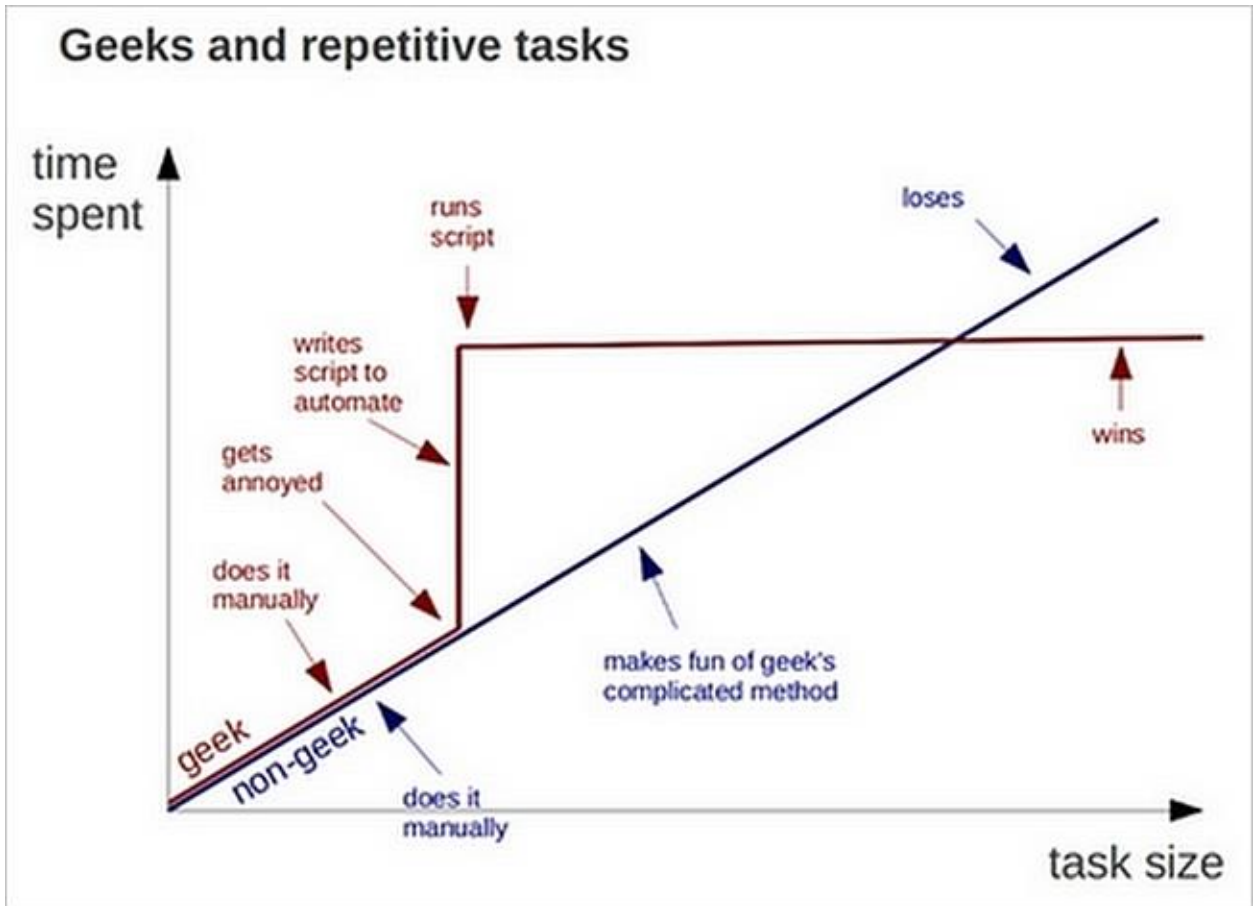8. If all goes well in step 6, we get this:



Notice that the 40 features have been placed in the feature tables (features (subsets)). The counter will also reflect the number 40. If you are happy at this point, you can proceed to click on "Export" and save the file in the ARFF format.



This is probably the most labor-intensive part of the homework because of the manual interaction required with LightSIDE. If you are inclined, you may use something like "Macro Express" to record your mouse & keyboard movements and play them back repeatedly on the different experiments. One key thing to remember when doing that is to have a consistent

window size when recording and playing back the macros. Another property of the lightSIDE GUI is that the number of rows shown to you on the features table is controlled by your window size. I won't be responsible for supporting your macros, but if you are familiar with them you should use it since the homework will have you do these tasks a fair number of times.
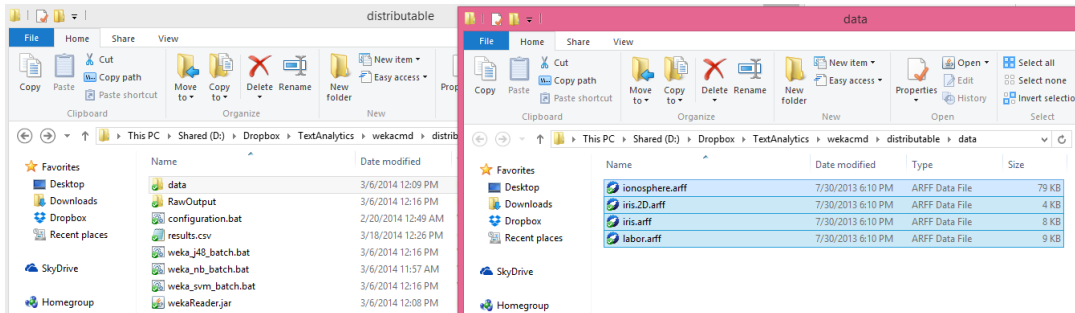


**Geeks and repetitive tasks**

9.  WEKA and DNDW
    Once you have installed Weka and DNDW, you can start processing the ARFF files with DNDW. DNDW will be recording your results in results.csv. You will be required to turn results.csv in along with your homework.
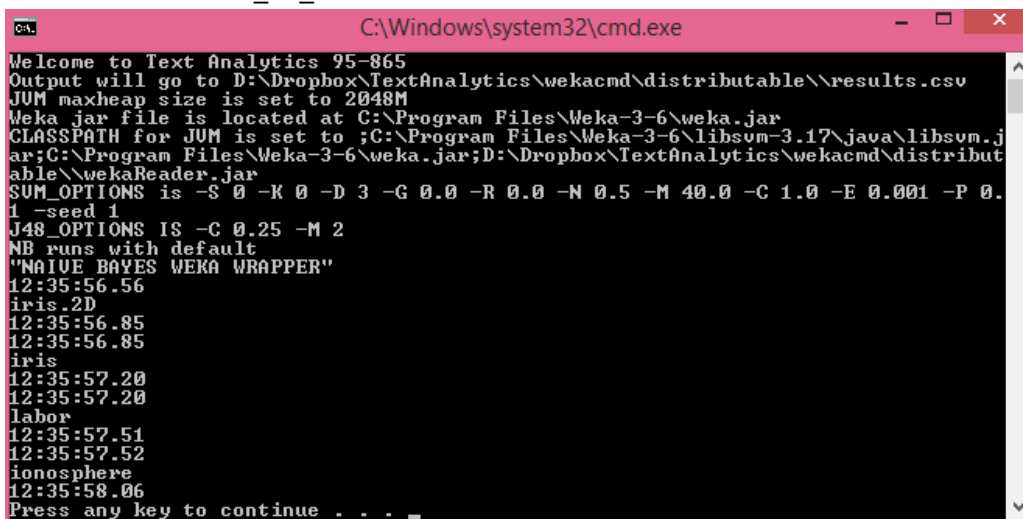
    Mac users will have to install Weka and DNDW on the Heinz virtual lab and process them there. Keep in mind that when using the Heinz virtual lab, your files will be lost after you disconnect from the machine. To avoid any issues relating to that, please put your files on a thumbdrive and connect it to the virtual lab machine.

10. DNDW – Drag and Drop Weka

Select your ARFF files, proceed to drag and drop them on the batch file that contains the algorithm you want to run.

For example, if we want to process all these ARFF files with the Naïve Bayes algorithm, you drop them on to the weka_nb_batch.bat file:



All the batch files will execute configuration.bat first and display the environment variables in the beginning.

The remaining outputs are the dataset name, start time and end time. All the raw outputs from Weka will be placed inside the ./RawOutput folder where you can examine them.

The program will append the results of your experiment to **results.csv**. You can open this file with excel later to collect your results (precision, recall, f number). This is done for your convenience, otherwise you would have to go through the Weka output and manually copy the values you need for your homework report. Keep in mind that you should not lock the file for editing when the batch files are running, otherwise your results won't be appended to the csv file. If you want to peak inside results.csv, making a copy would be the best way.

**Notes for Mac Users**

Weka, LightSIDE and libSVM are java based programs. This is good news because they can be run natively on your computer. The automation component (DNDW) is unfortunately written in windows batch files. You have a few options:

1. Use Bootcamp and run everything in windows.
2. Use OSX Lightside and create the ARFF files in Mac, copy them into Bootcamp and run DNDW in windows bootcamp.
3. Use OSX LightSIDE and create ARFF files in mac, copy them on to a USB drive and run DNDW in Heinz Virtual Labs.

Instructions for Heinz Virtual Labs can be found here:

http://www.heinz.cmu.edu/computing-services/virtual-labs/index.aspx