GETTING STARTED with the IMDB data
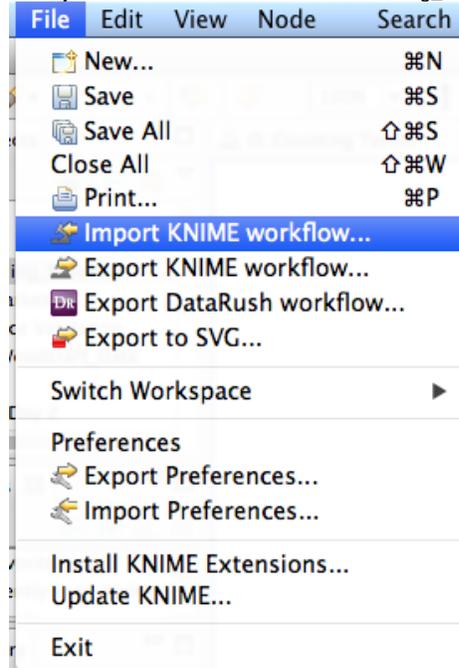
1) download the IMDB data located at http://abbottanalytics.com/data/textminingdata/IMDB%20Review%20Data.zip

UNZIP the data! In the class, we will be using the POS and NEG reviews with 100 reviews in each directory

2) download the baseline workflow, http://abbottanalytics.com/data/textminingdata/IntroTextMining_Baseline.zip
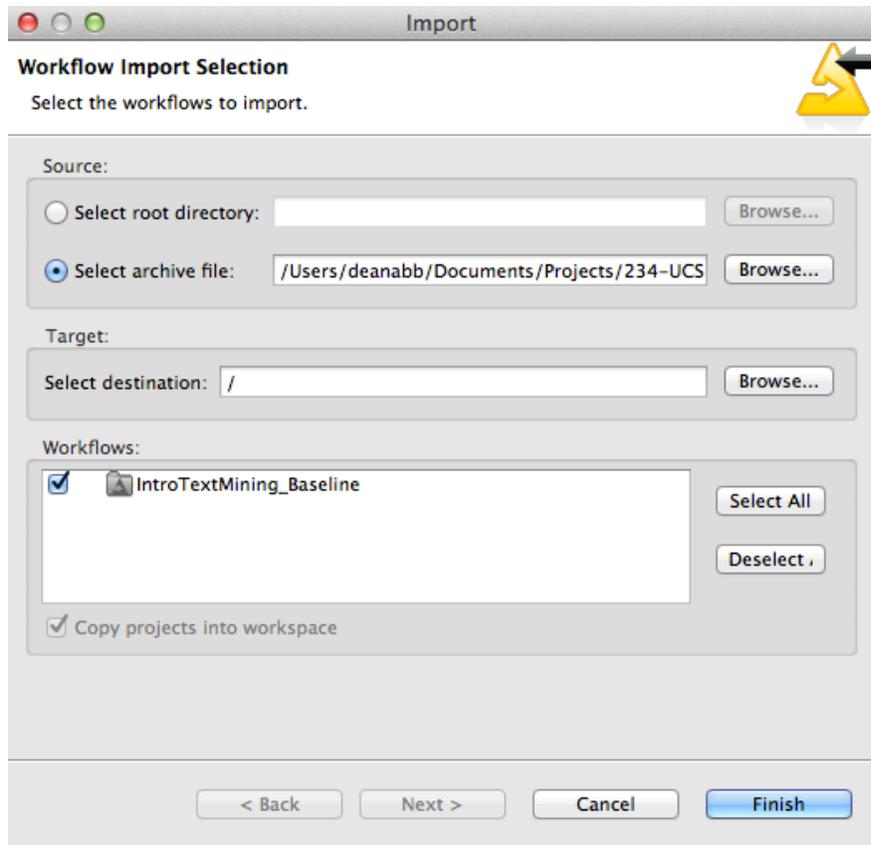
DO NOT UNZIP THIS FILE. You will be importing it as an **archive**

Once you have downloaded the IntroTextMining_Baseline.zip workflow, open the KNIME program, and import the KNIME workflow
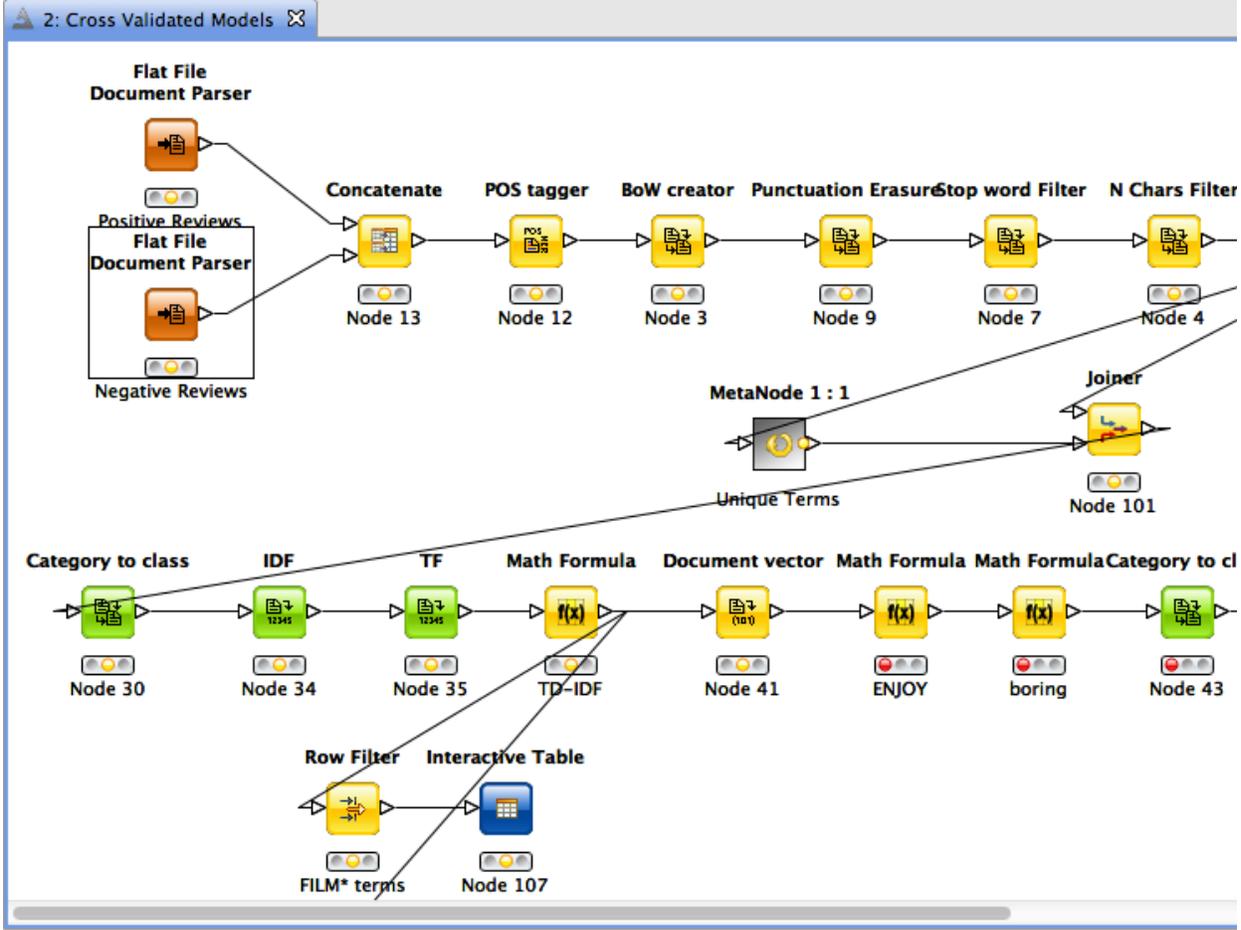


After clicking on Import KNIME workflow…, you will see a window that allows you to select the workflow to import. Click on the "Select archive file:" radio button, find the workflow in your file system, and by default, you should see the "IntroTextMining_Baseline" workflow checked. Click "Finish". Leave the "Select destination" directly alone--this just pushes the workflow to your default workspace.

If you UNZIPped the workflow, you can import it as a "root directory" by pointing to the top level directly containing the unzipped workflow.
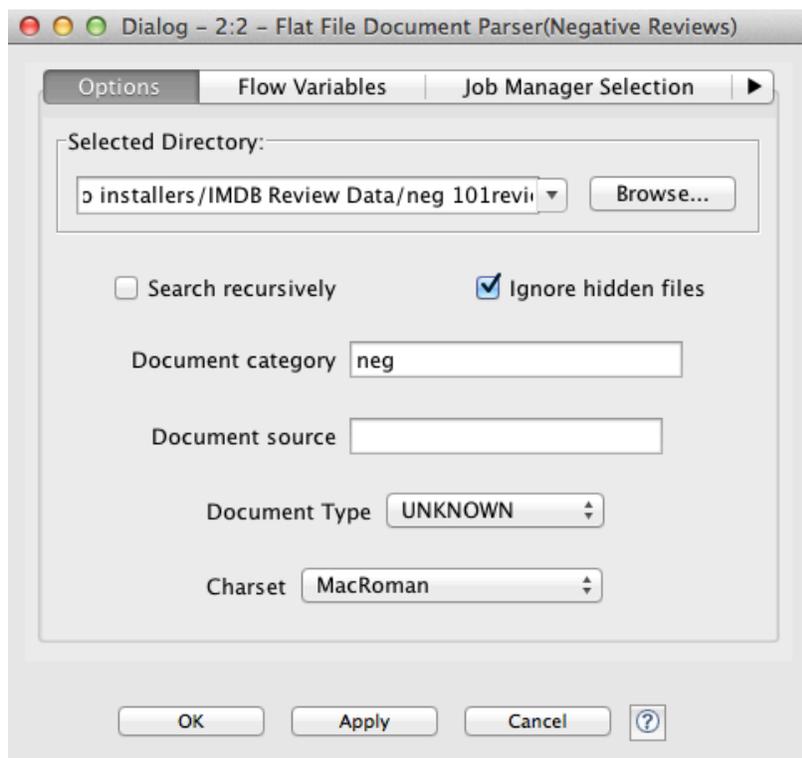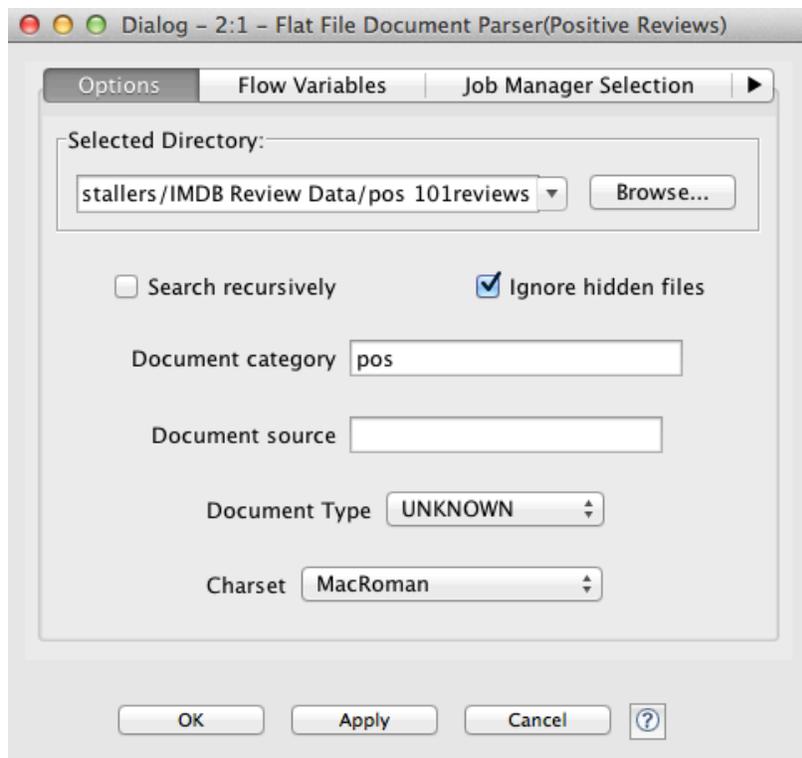
You will see the workflow in the left side of your main KNIME window in the "Workflow Proejcts" area. Double-click on the workflow. It will load and look something like this:

**2: Cross Validated Models**

Flat File Document Parser — Positive Reviews

Flat File Document Parser — Negative Reviews

Concatenate — Node 13
POS tagger — Node 12
BoW creator — Node 3
Punctuation Erasure — Node 9
Stop word Filter — Node 7
N Chars Filter — Node 4

MetaNode 1 : 1 — Unique Terms

Joiner — Node 101

Category to class — Node 30
IDF — Node 34
TF — Node 35
Math Formula — TD-IDF
Document vector — Node 41
Math Formula — ENJOY
Math Formula — boring
Category to cl — Node 43

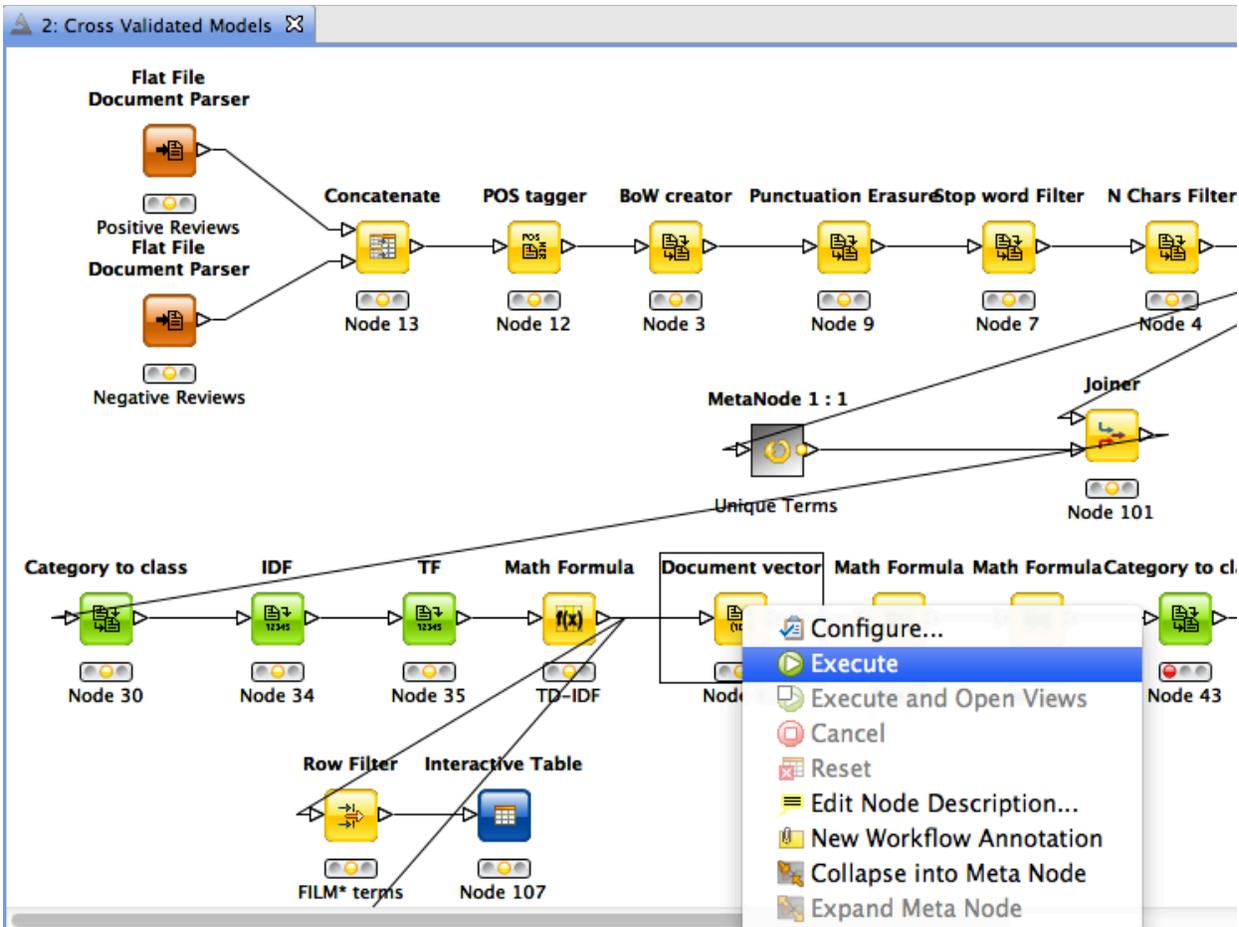Row Filter — FILM* terms
Interactive Table — Node 107

If you haven't installed the extensions for TextProcessing and Math Expressions, you will see some icons that are generic. These will not execute. Please go back to the installation document and install these two extensions.
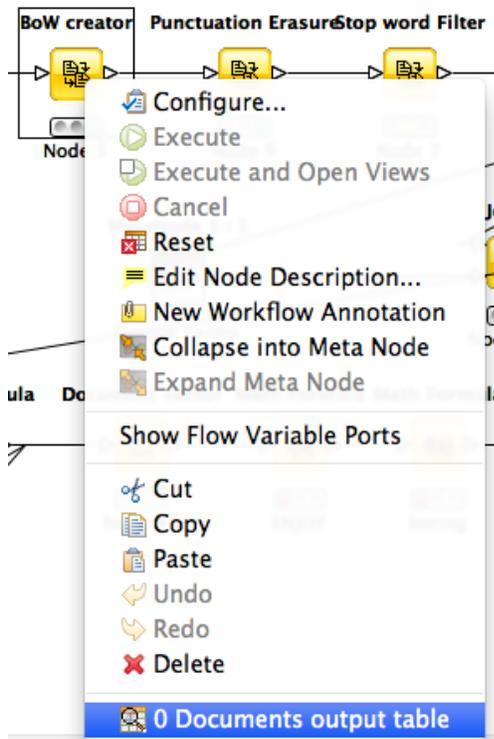
The configuration issue is pointing to the right data. Open the icon that says "Positive Reviews". Set the path to the data to point to the IMDB folder "pos 101reviews". Also make sure you set the Document category as "Pos". This merely labels the positive reviews with the label "Pos", indicating these are positive review. Do the same for the negative reviews in folder "neg 101reviews". Click OK.

Dialog – 2:1 – Flat File Document Parser(Positive Reviews)

Options    Flow Variables    Job Manager Selection    ▶

Selected Directory:

stallers/IMDB Review Data/pos 101reviews  ▼    Browse...

☐ Search recursively        ☑ Ignore hidden files

Document category  pos

Document source

Document Type  UNKNOWN  ⬍

Charset  MacRoman  ⬍

OK        Apply        Cancel    ?



Dialog – 2:2 – Flat File Document Parser(Negative Reviews)

Options    Flow Variables    Job Manager Selection    ▶

Selected Directory:

ɔ installers/IMDB Review Data/neg 101revi  ▼    Browse...

☐ Search recursively        ☑ Ignore hidden files

Document category  neg

Document source

Document Type  UNKNOWN  ⬍

Charset  MacRoman  ⬍

OK        Apply        Cancel    ?

Once you do this, the node should have the yellow dot illuminated (it was red before). We will walk through each step in class, but to test the workflow, right-click on a node that is yellow, such as the "Document Vector" node and select "Execute". Each node as it completes will have its dot turn green.

At any time, you can still click on nodes. In fact, after a node is green, you can right-click on that node and select the bottom entry which is the table view of that node. For example, if you want to look at BoW creator results, you can right-click on that node and select the "output table" option.

In my data, the results look like this. The POS refers to the Part of Speech label. The first entry is an adjective, and that POS has been appended to the term itself.