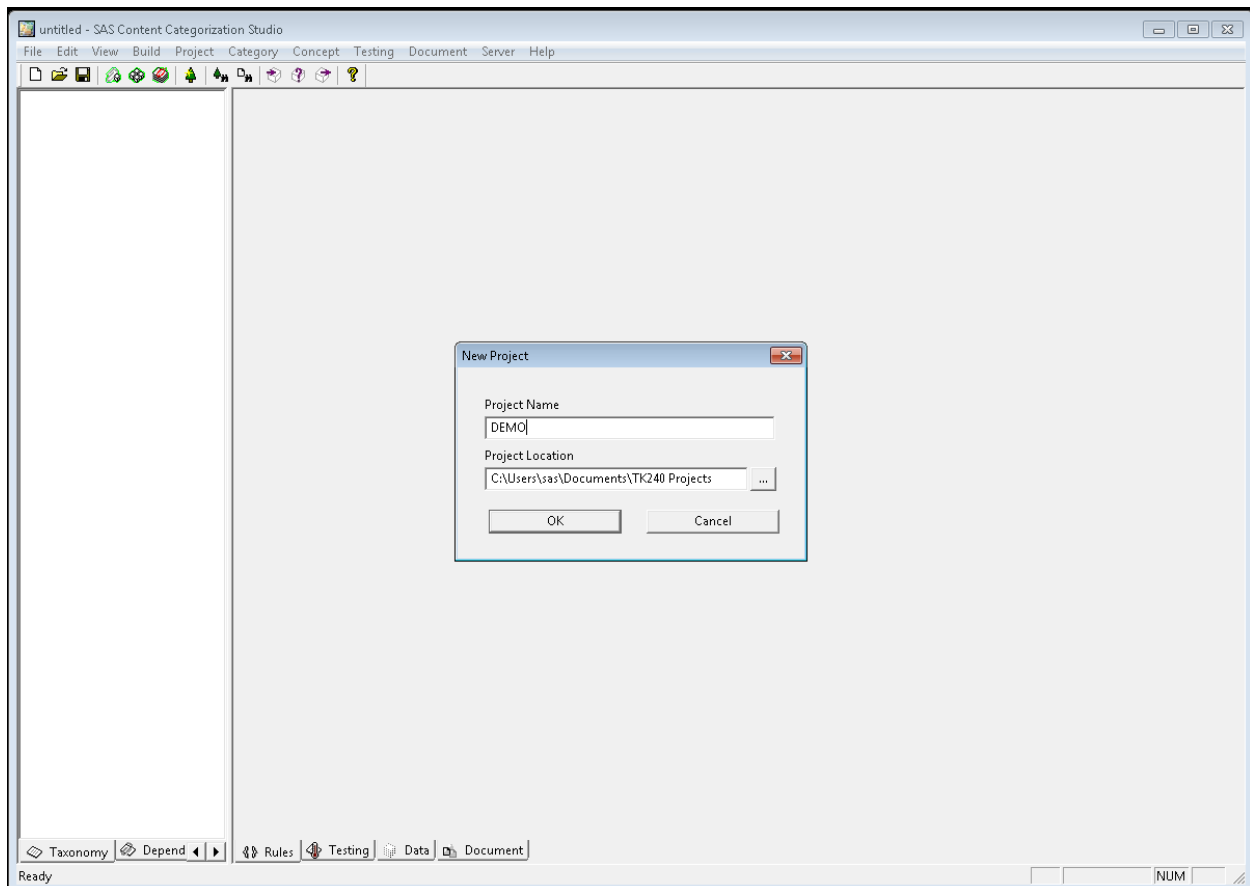Here are steps to create a SAS Content Categorization project.  There will be steps to create both a category rule as well as a concept rule.  Each rule will be built or compiled, and then testing on the rule will be done on a set of documents that is part of the zip file.
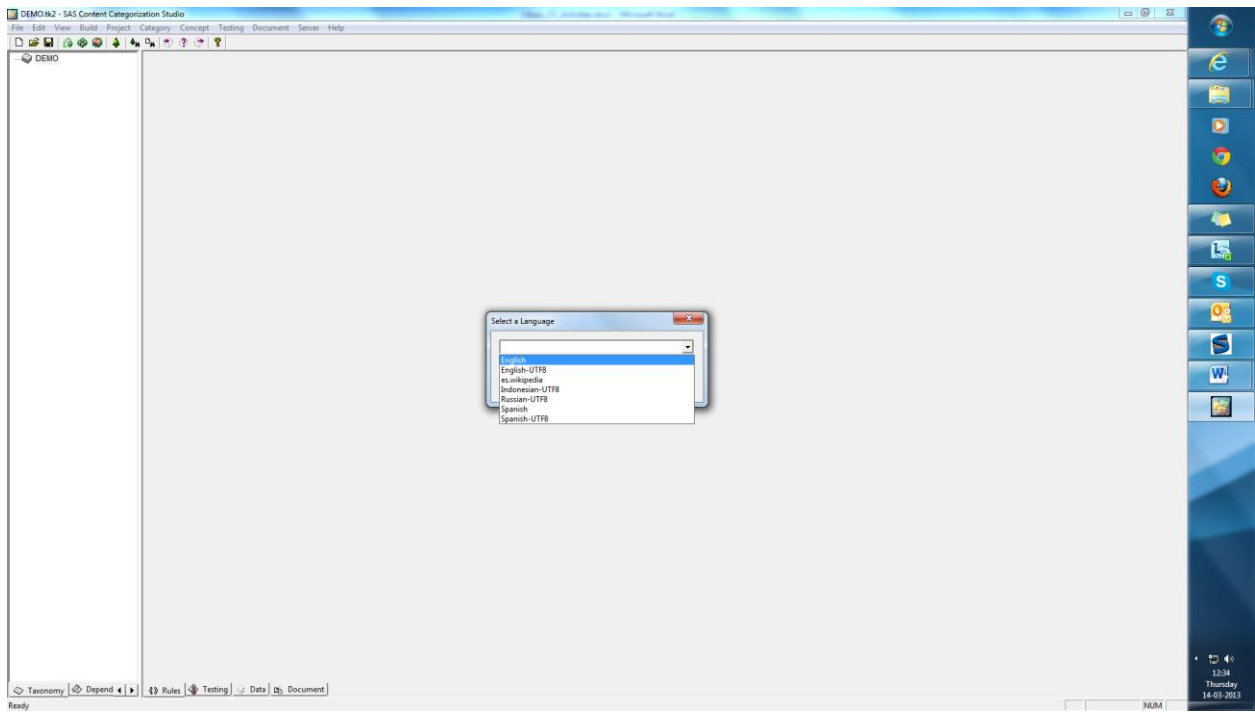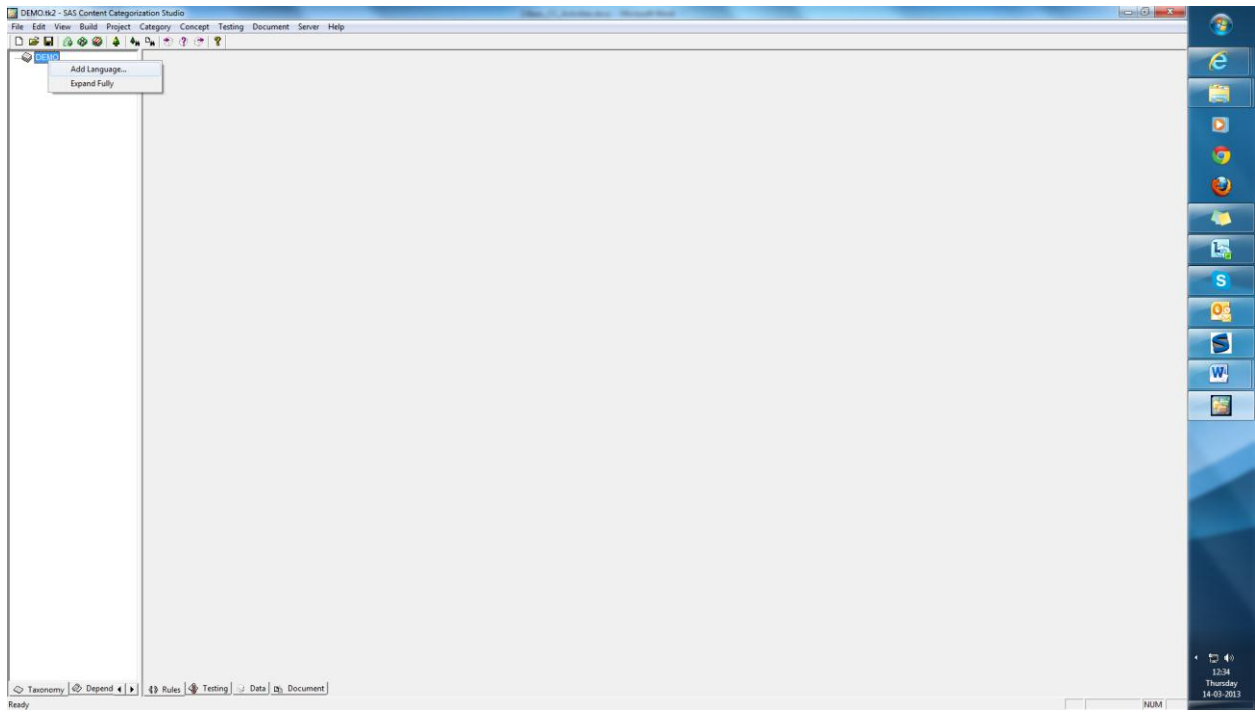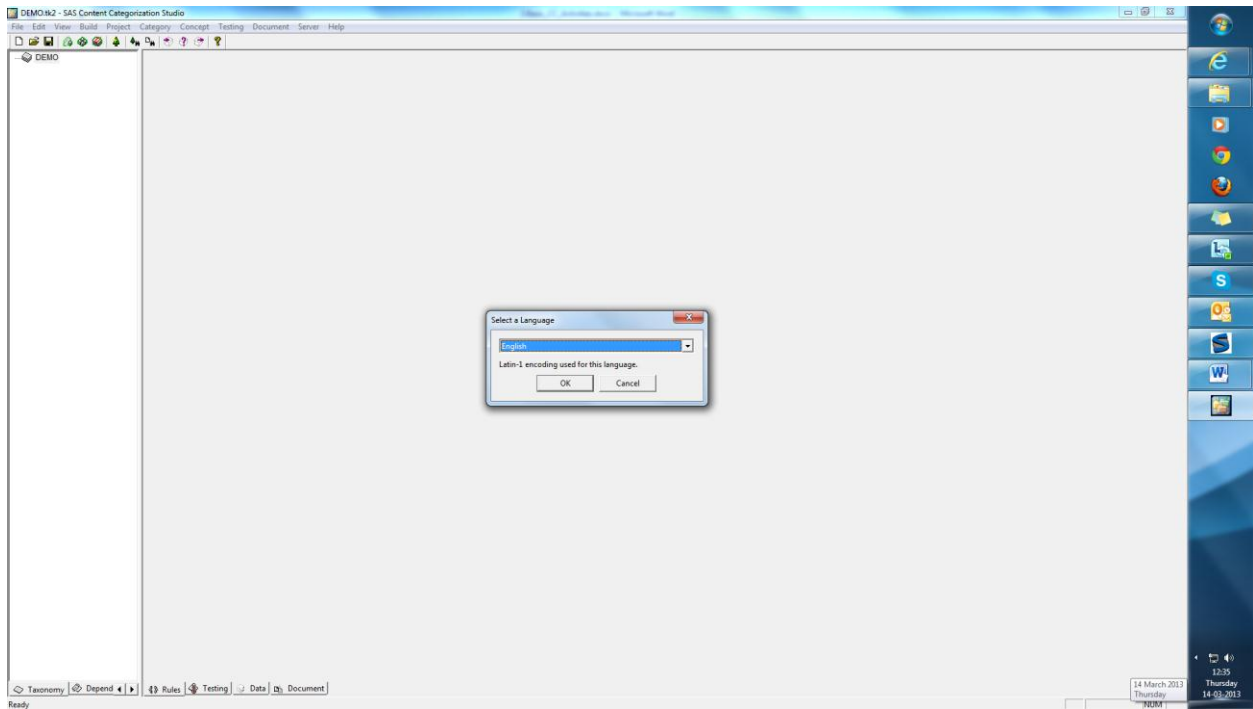
## Part 1:  Creating a project

1.  Create New Project and set the Location where the new Project Folder will be stored.



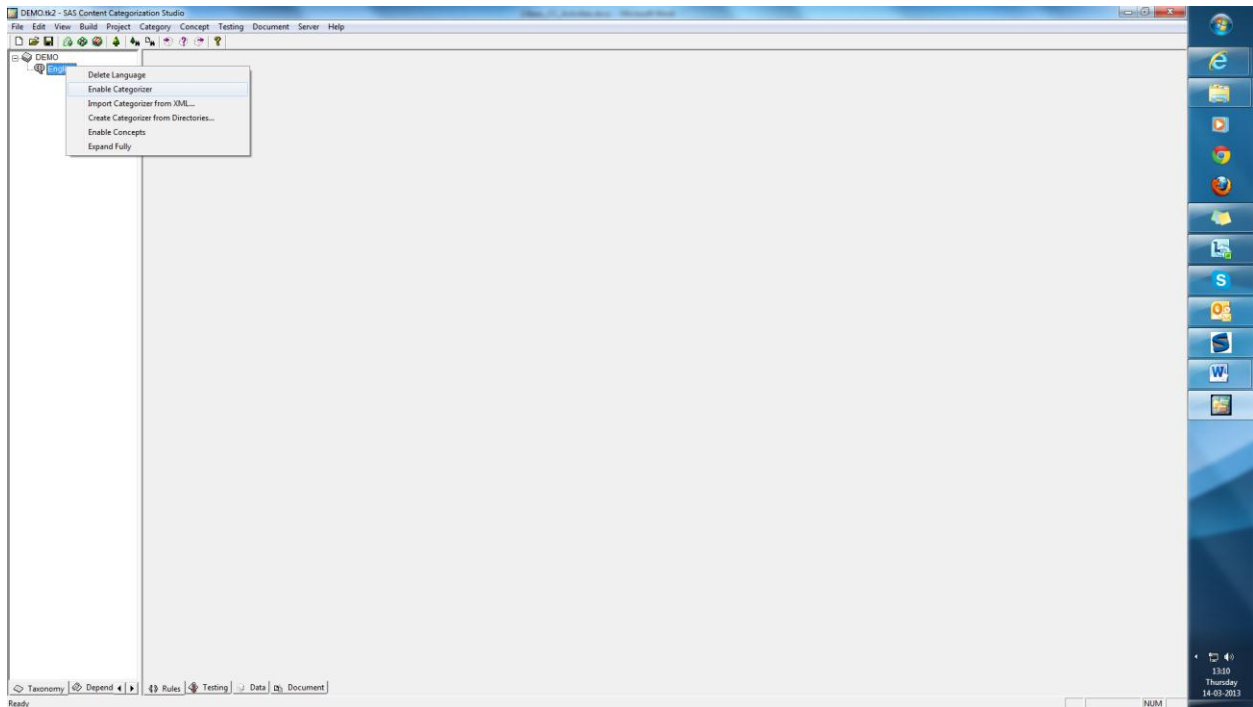Now we will begin creating rules in the newly created *Demo* project.

2.  Add language to the new project *Demo*

DEMO.tk2 - SAS Content Categorization Studio

File Edit View Build Project Category Concept Testing Document Server Help

DEMO

Add Language...
Expand Fully

Taxonomy | Depend | Rules | Testing | Data | Document
Ready | NUM

---

DEMO.tk2 - SAS Content Categorization Studio

File Edit View Build Project Category Concept Testing Document Server Help

DEMO

Select a Language

English
English-UTF8
es.wikipedia
Indonesian-UTF8
Russian-UTF8
Spanish
Spanish-UTF8

Taxonomy | Depend | Rules | Testing | Data | Document
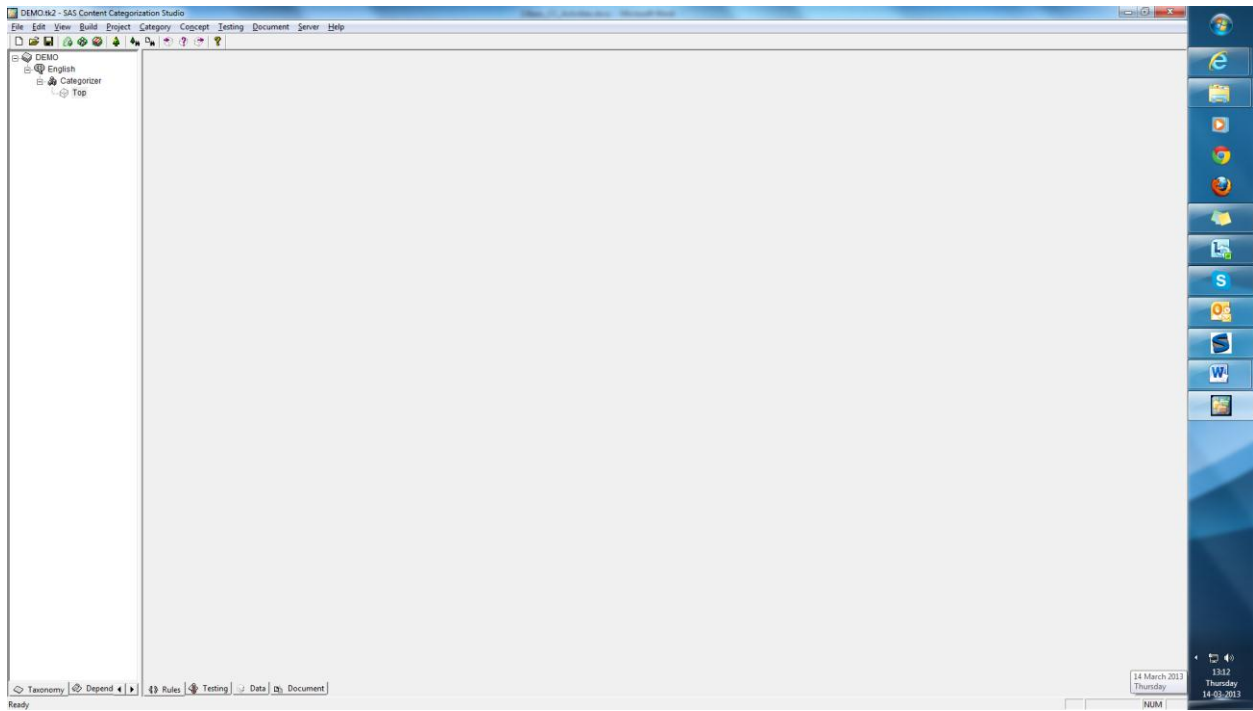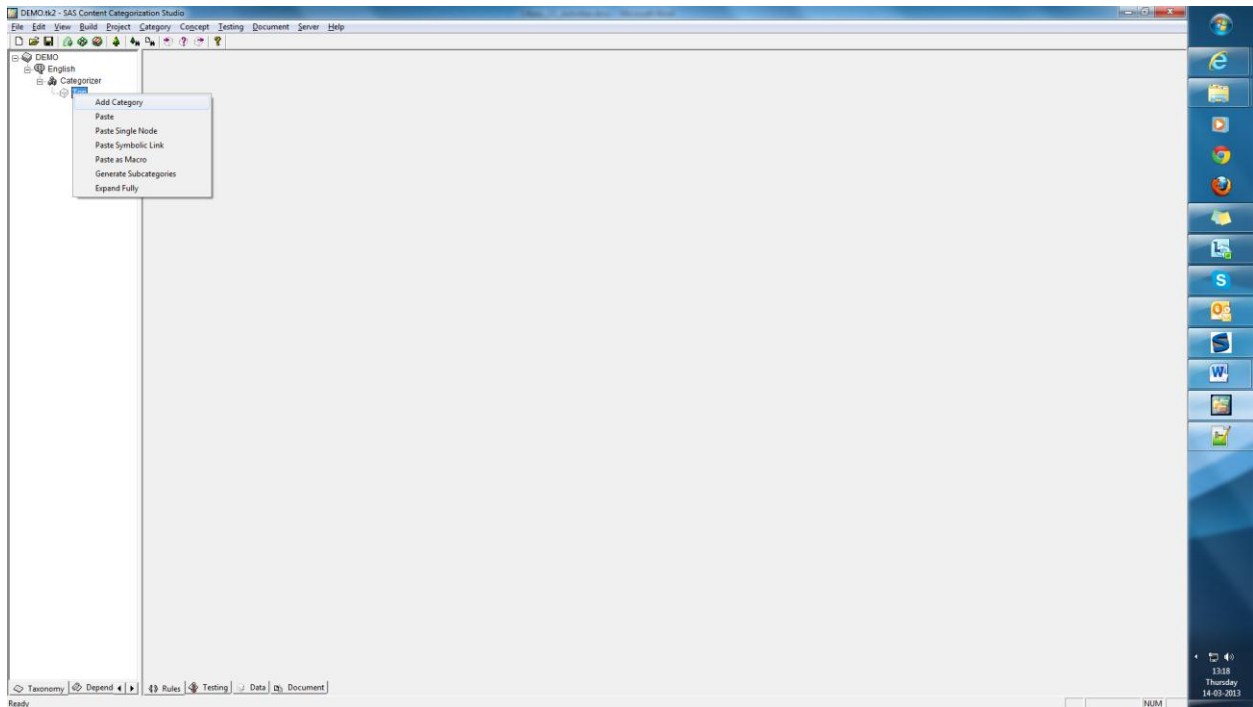Ready | NUM

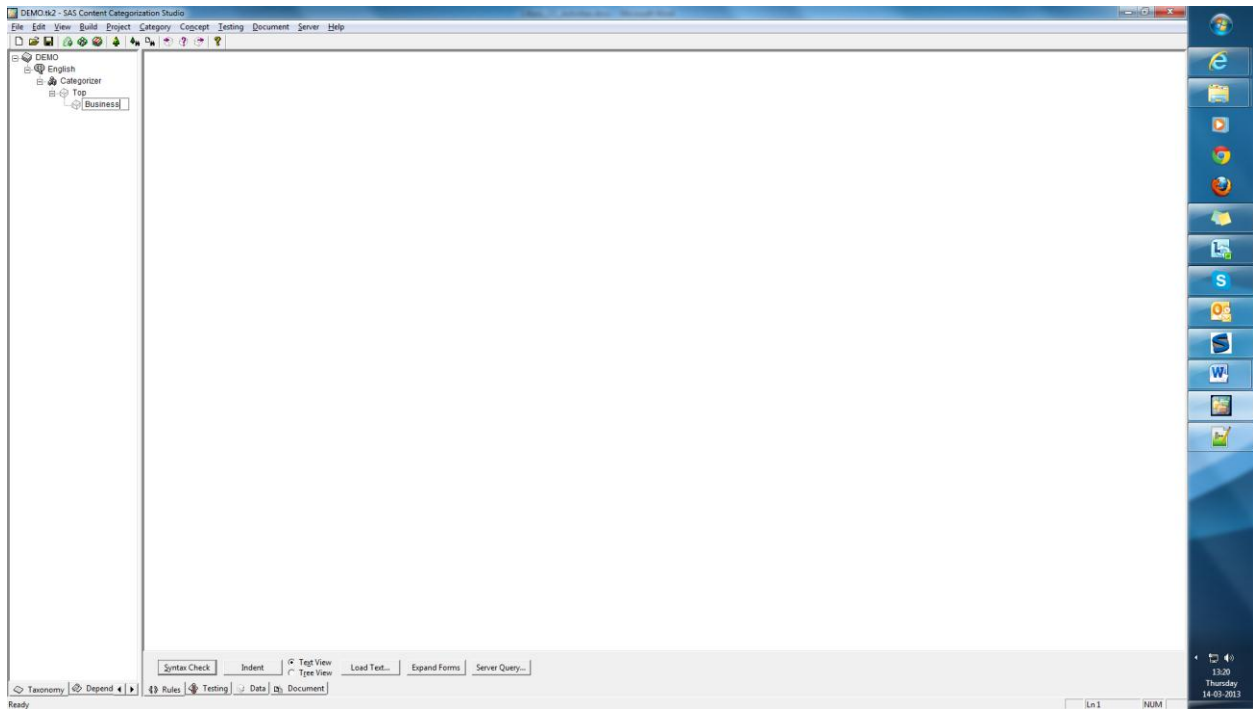## Part 2: Creating a Category

1. To create a category taxonomy, right click on English and choose Enable Categorizer.

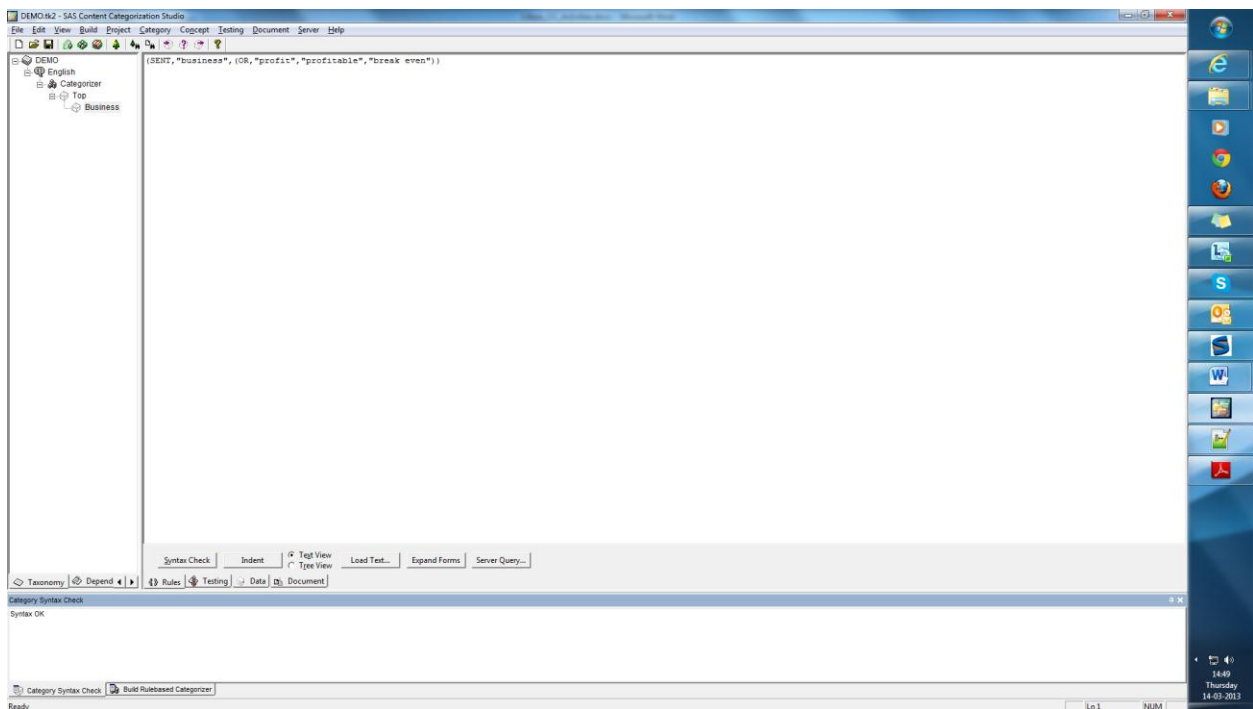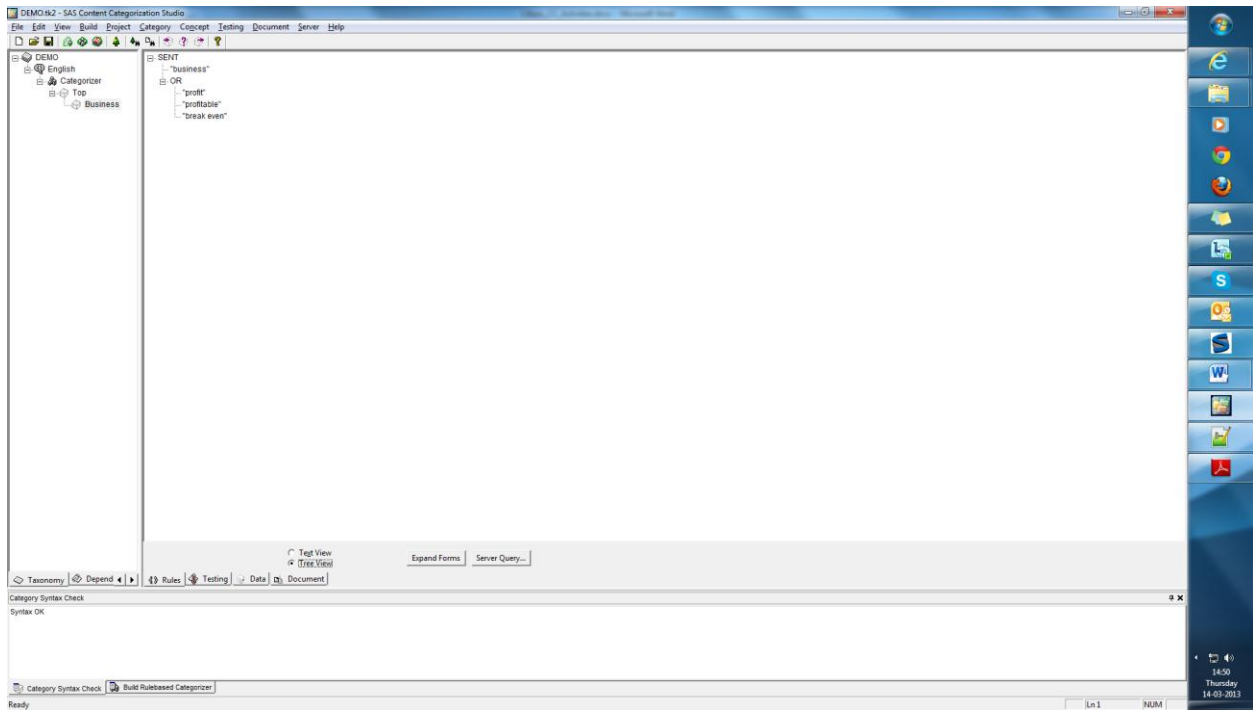2. Next, create a category by right mouse clicking on Top and select Add Category

On the rule tab of the Business category write the rules using Text View or Tree View.

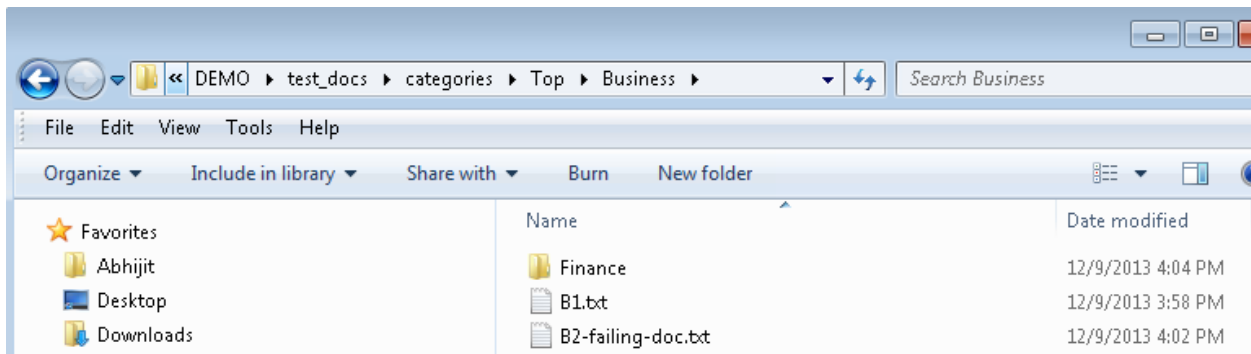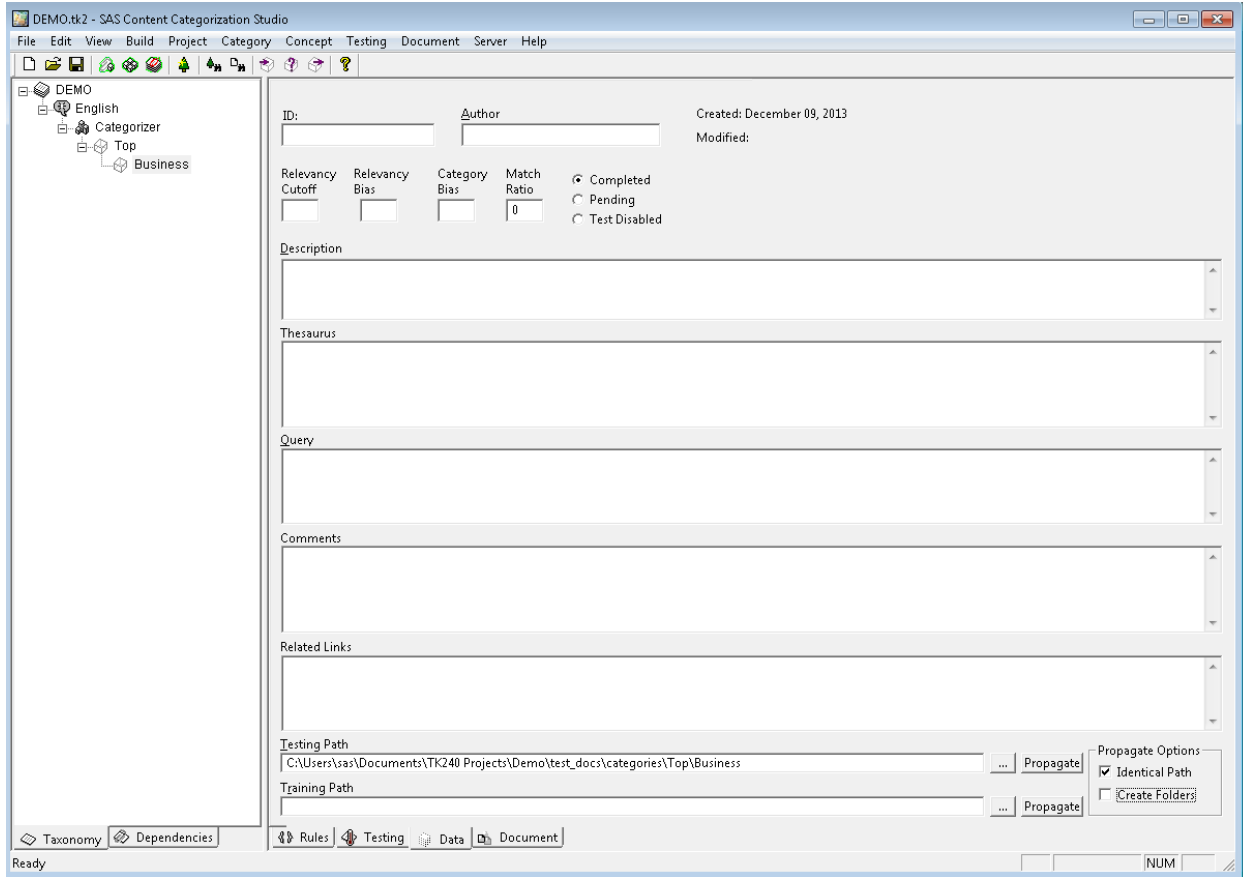3. Within Text View, in the rule tab insert the following rule:

```
(SENT,business,(OR,profit,profitable,break even))
```
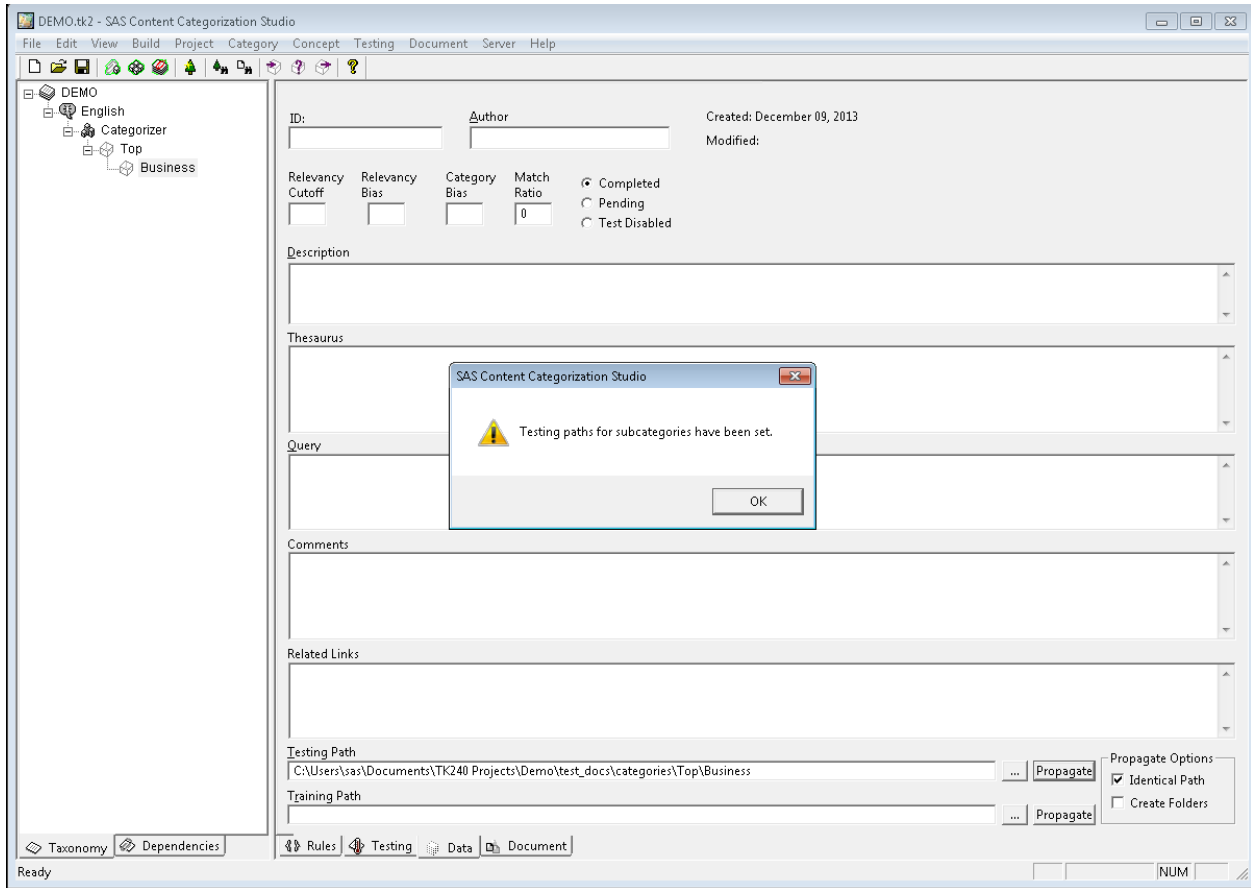


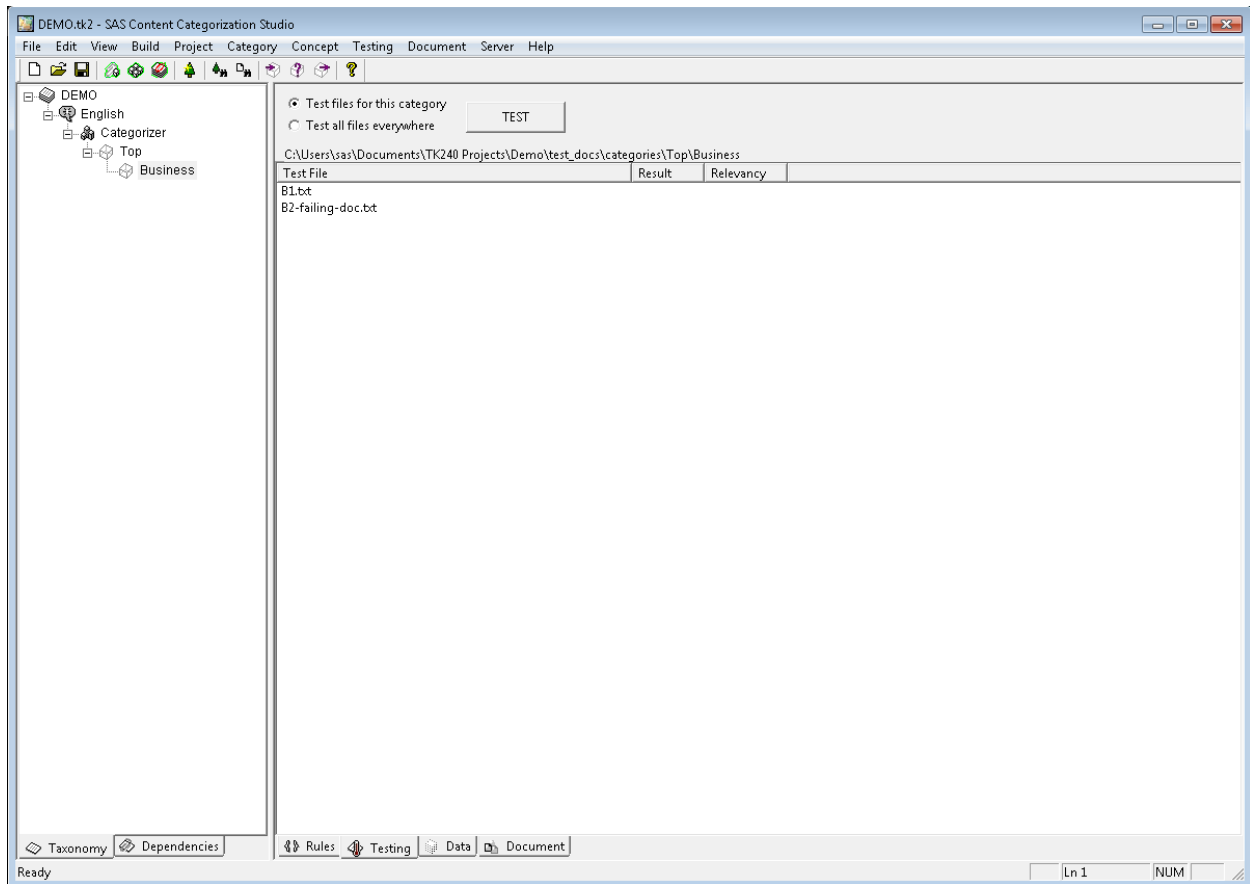Below is a screen shot of the tree view of the same rule.

4. On the data tab, set the path of the folder where the files to be tested are located and then click Propagate.
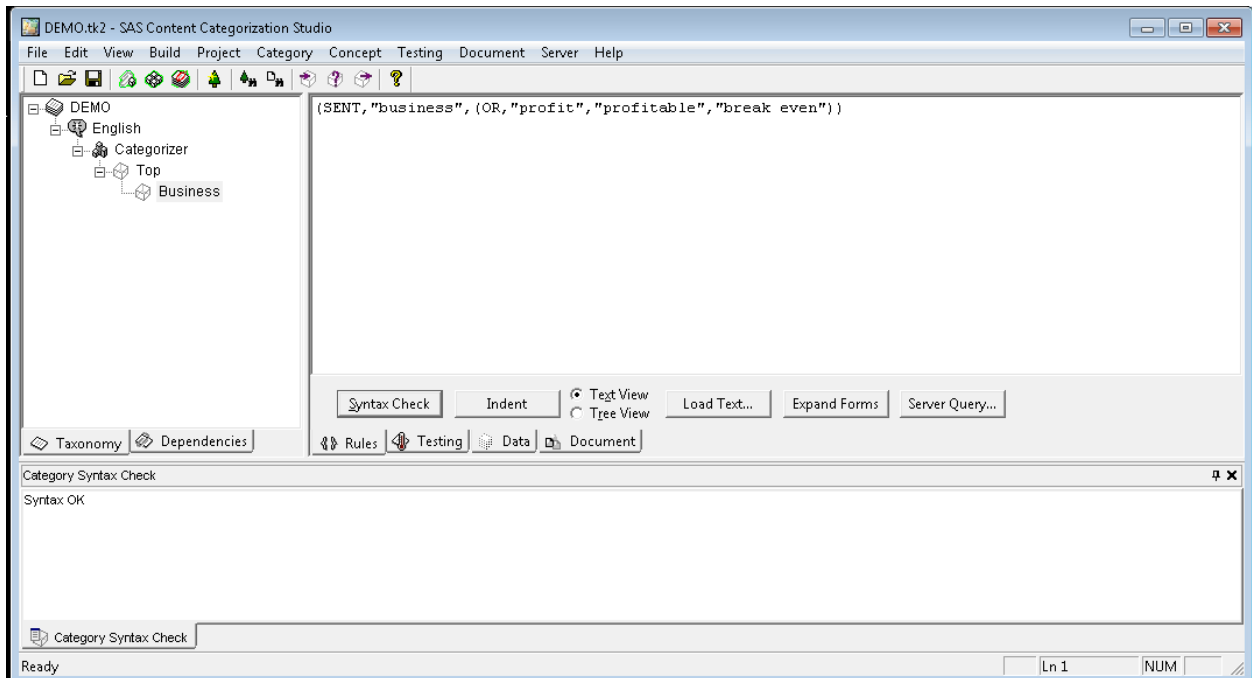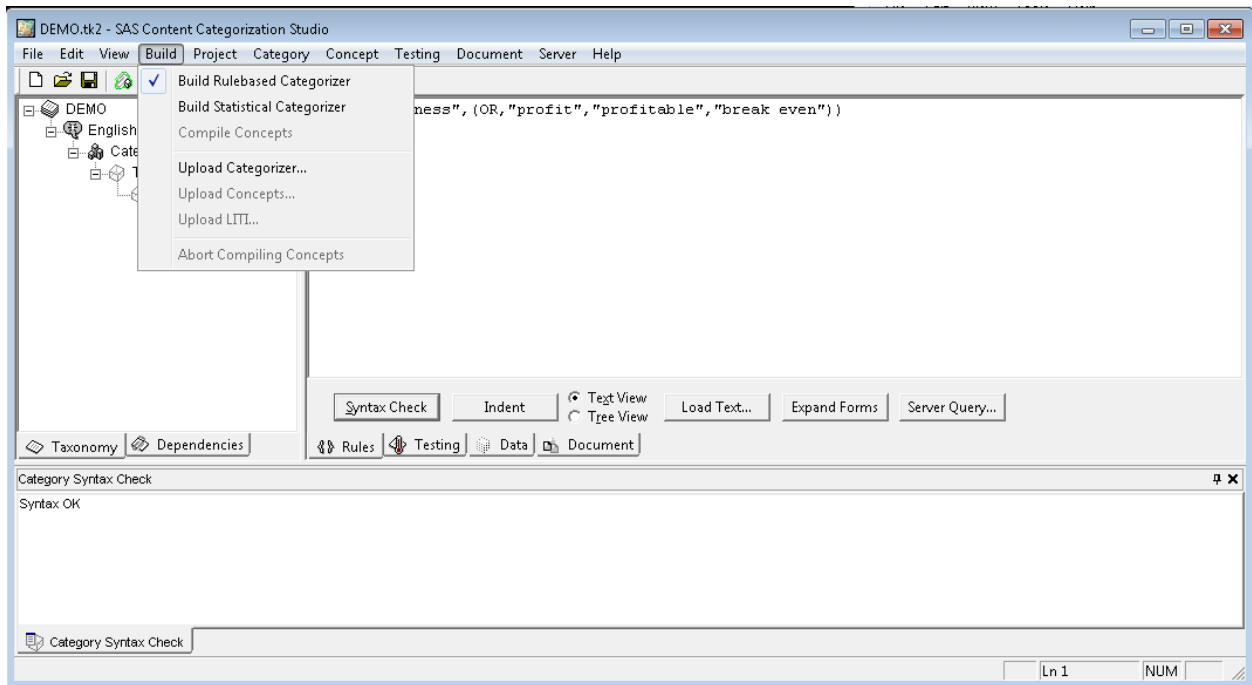
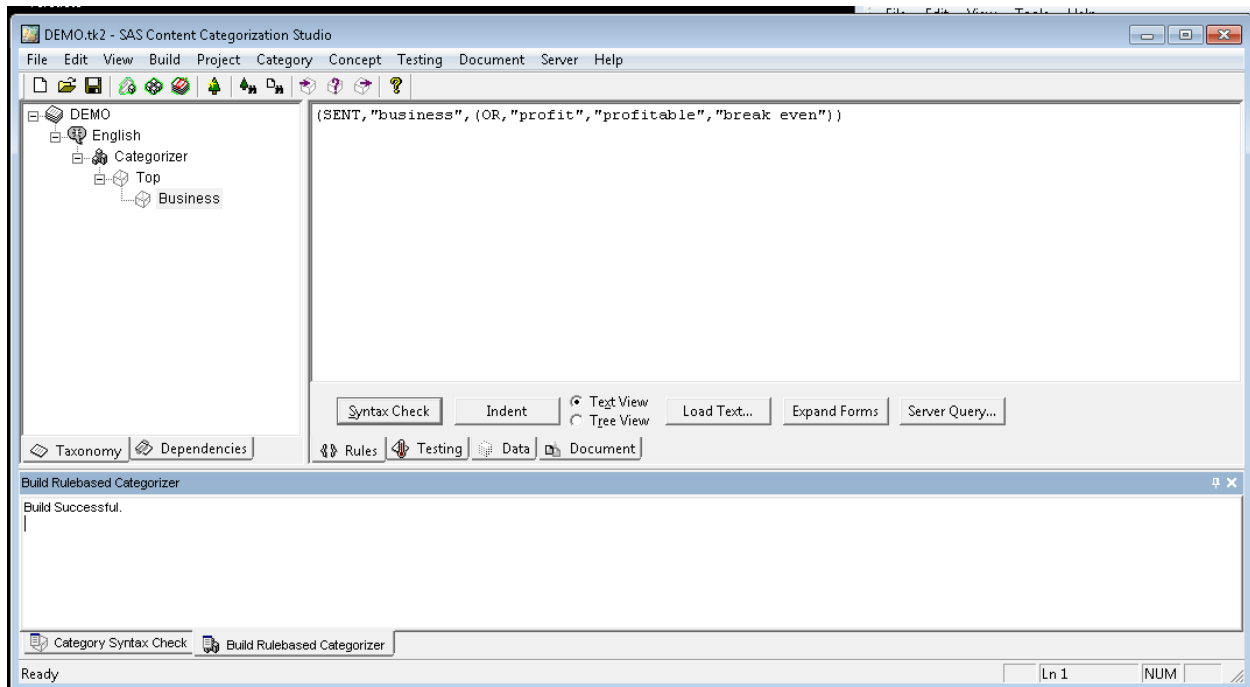The files in the testing folder will now be visible in the testing tab.

5. Now, build the category rules we created. To do that, check for the syntax of the rules using Syntax Check before building categories.
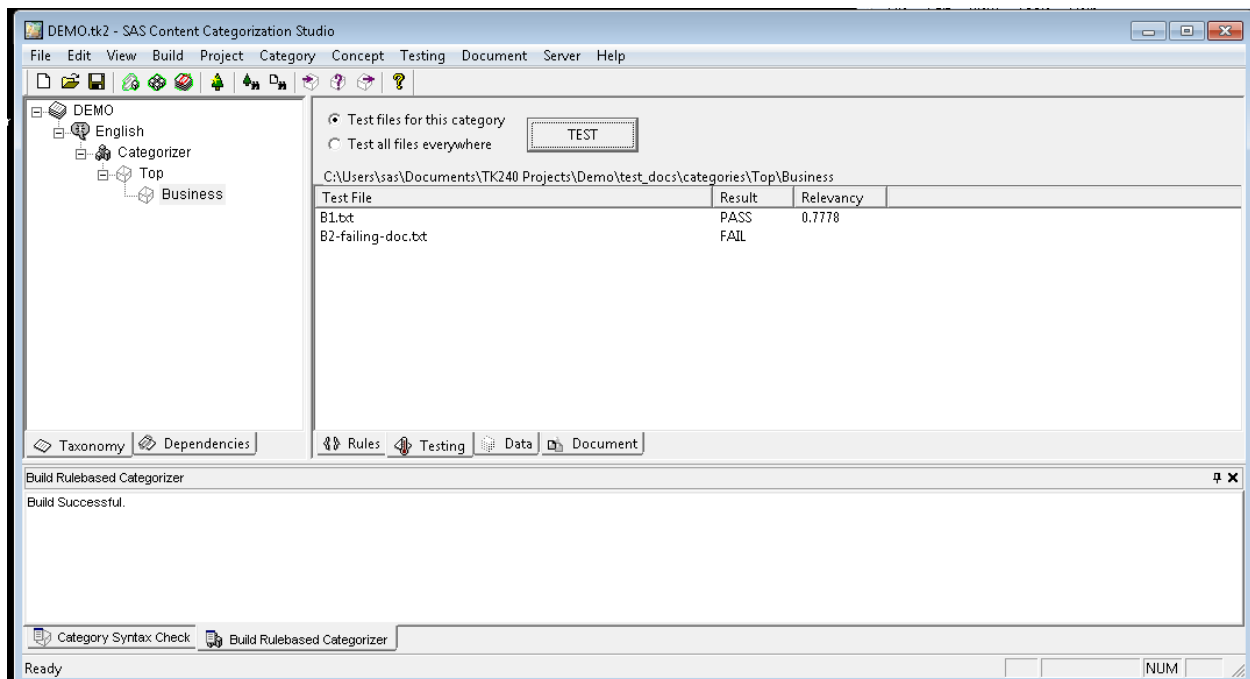
6. Above we created a boolean rule based category called Business, hence we will build the category using Build->Build Rulebased categorizer from the menu bar.
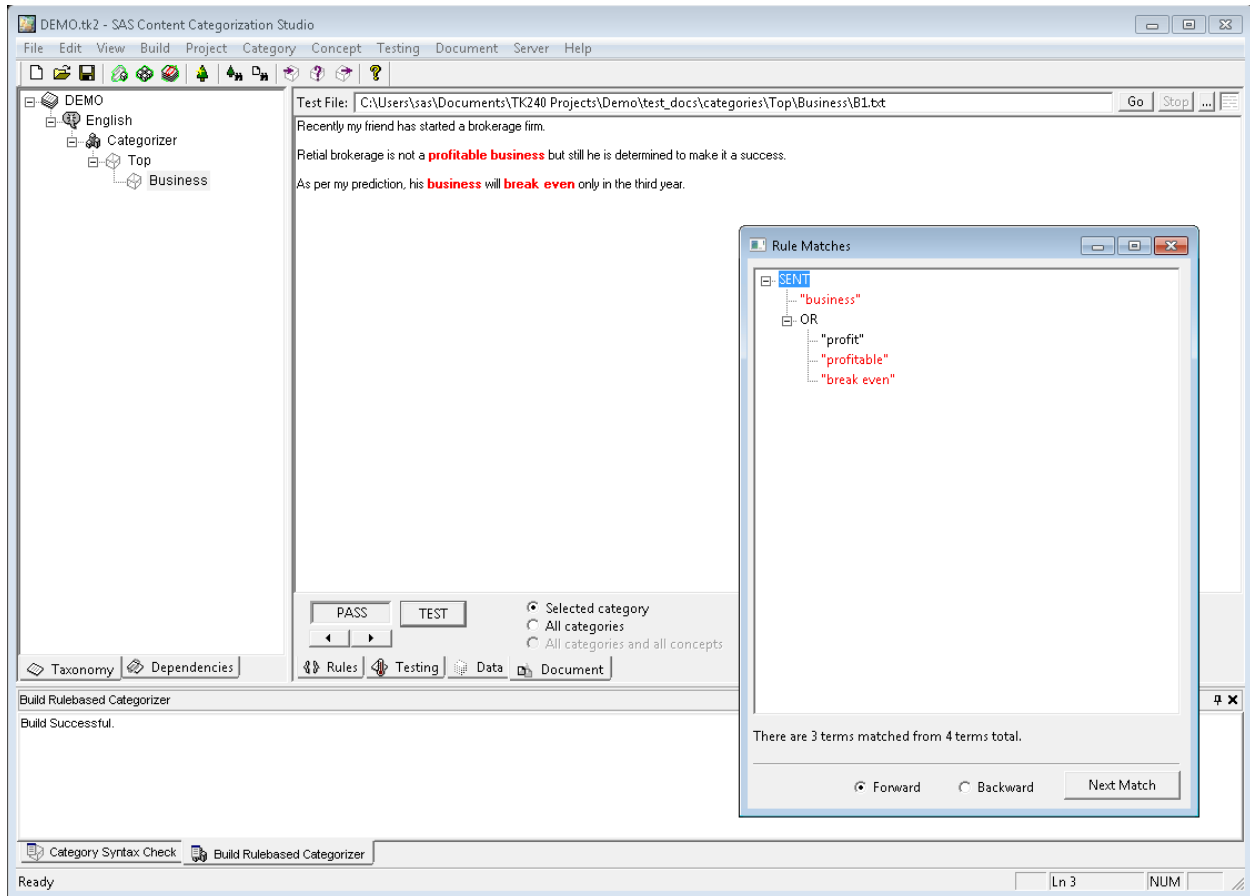


Build is successful.

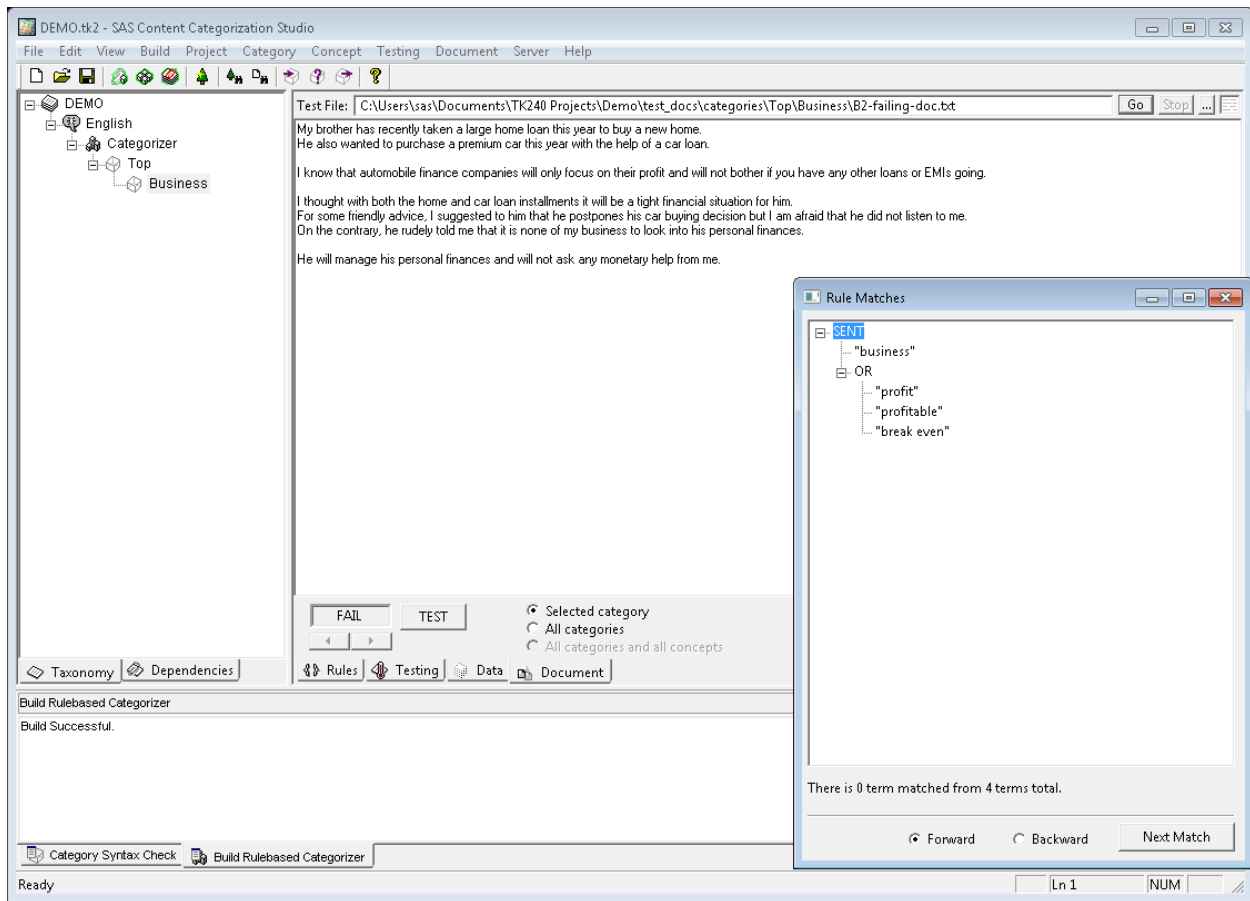7. Now to test the files go to the Testing tab and click the Test button.



We can double click any of the test documents shown above to see which words get matched or if the test documents did not match. Matched words will get highlighted in red.

Relevancy score shown above is useful to determine best match when the document matches more than one category.
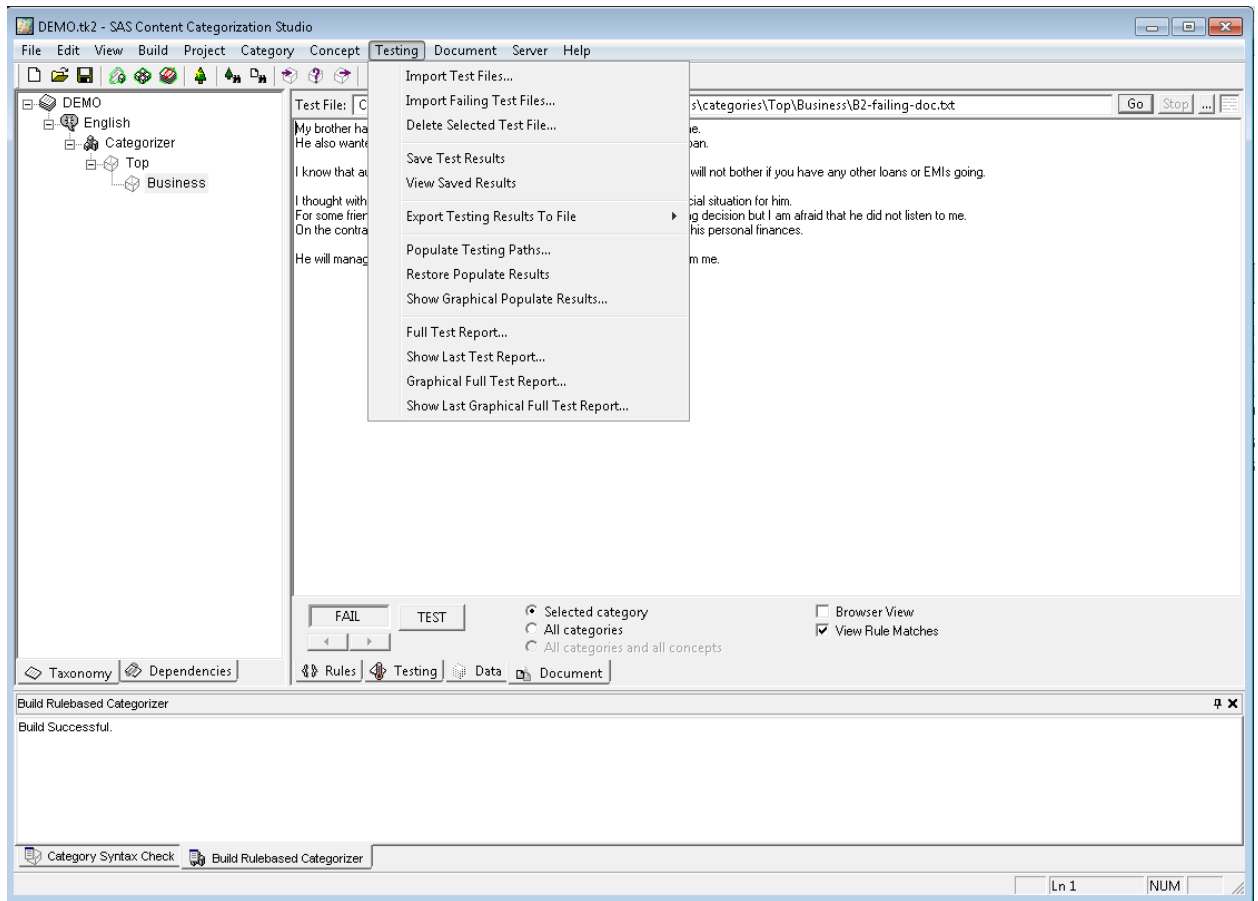
8. Double click the B1.txt to open that file on Document tab and check which words get highlighted.
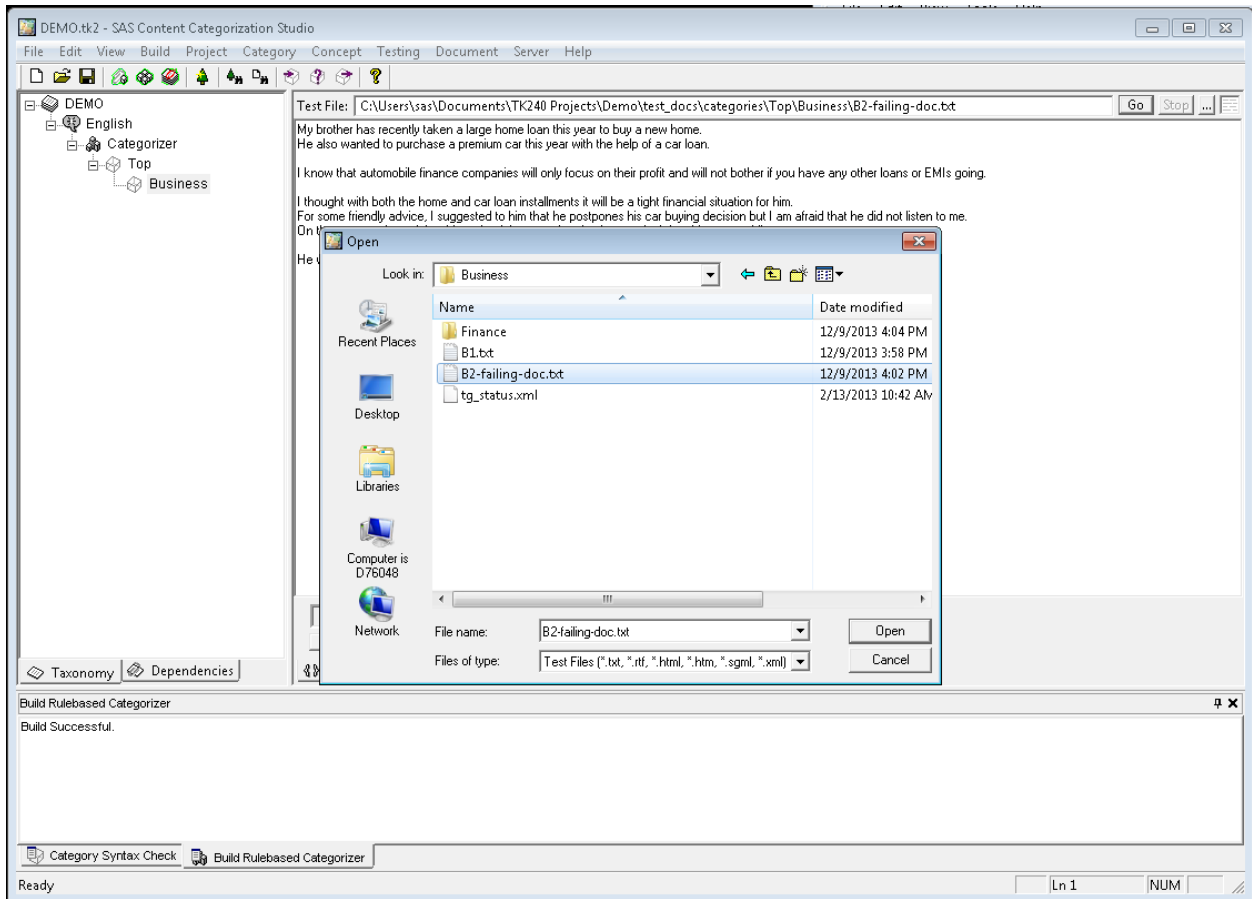


The file *B2-failing-doc.txt* fails because the terms profit and business do not match the rule defined.
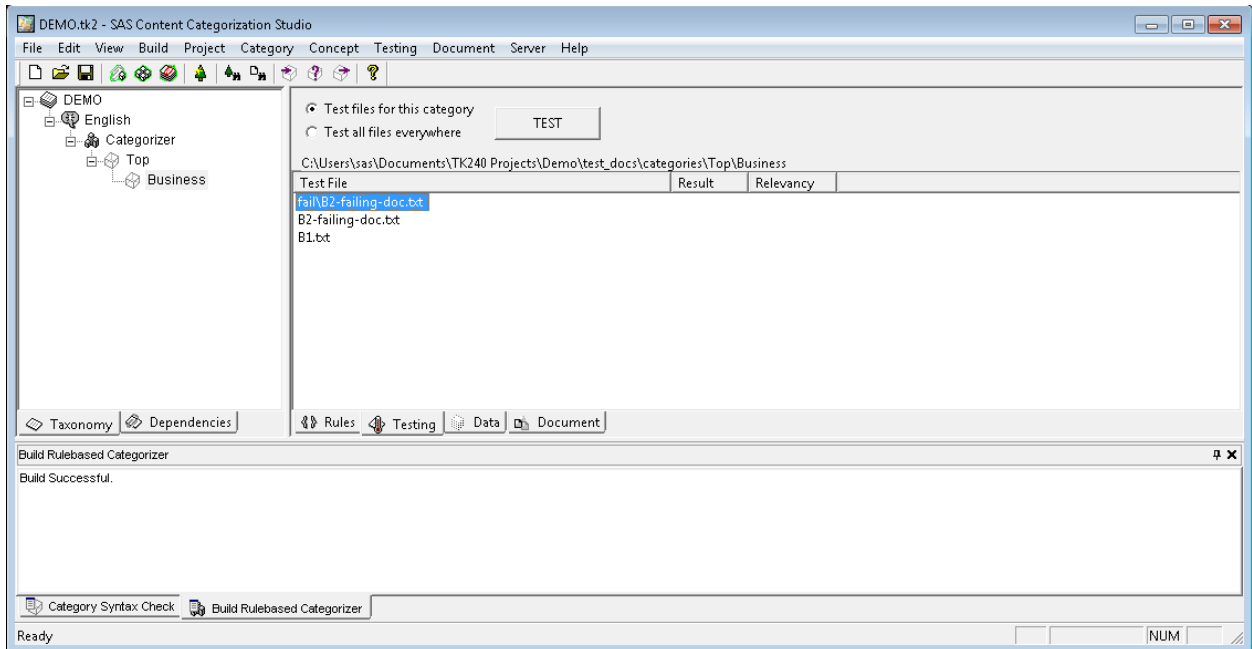
9. A good model should match all the relevant test documents (recall) while not matching irrelevant documents (precision). To test the failure of such irrelevant documents we can also check using Testing->Import Failing test Files….
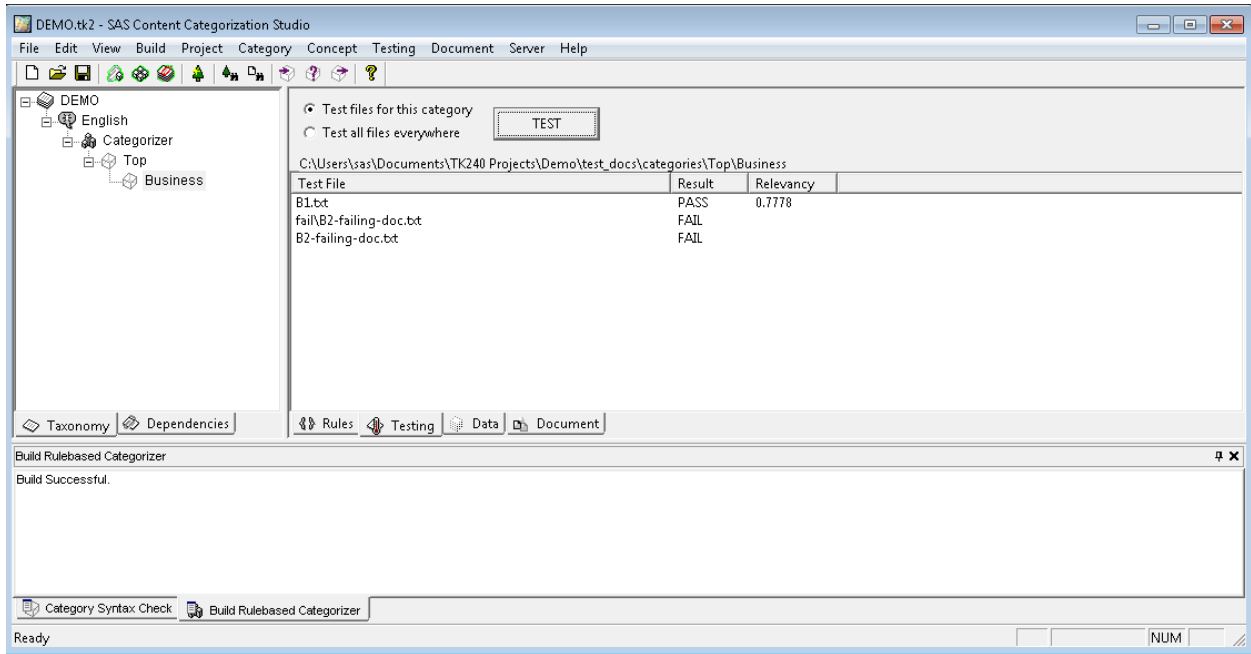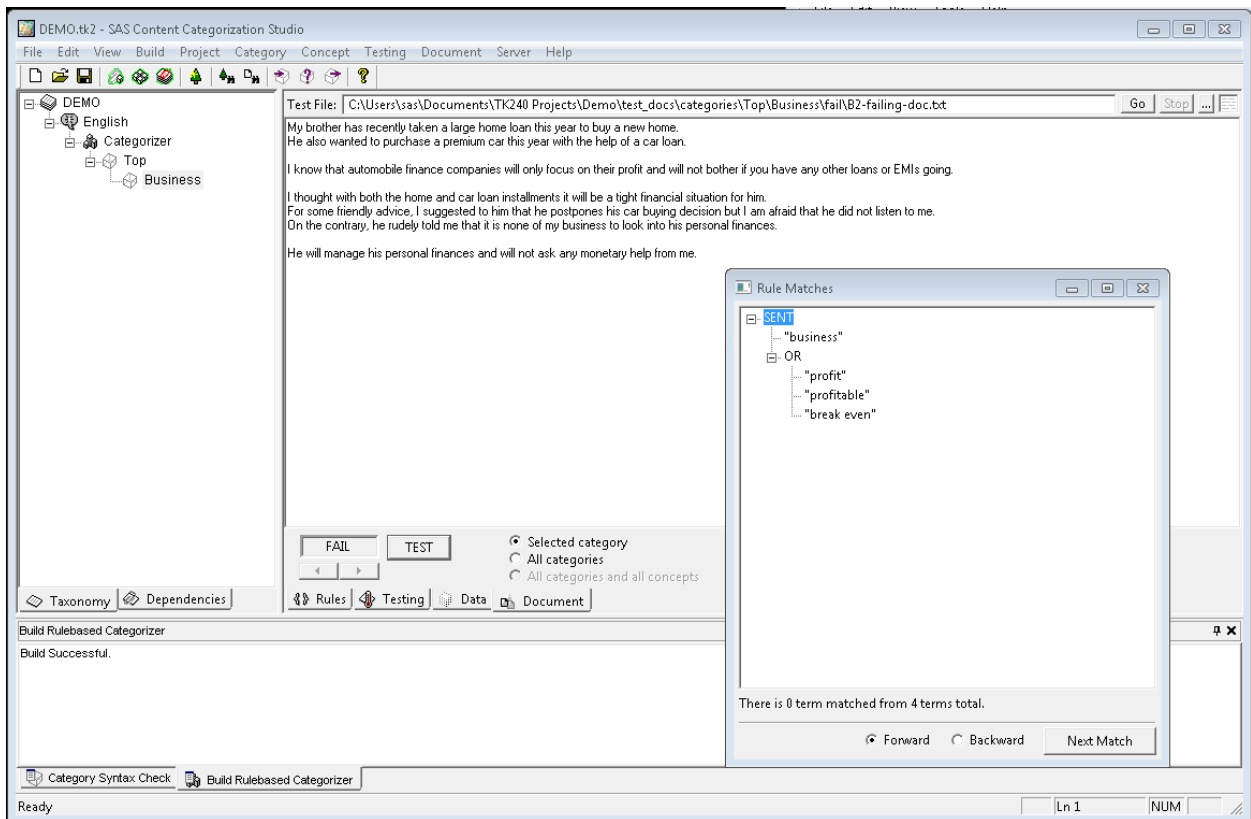
10. Open the Failing doc to be tested

The imported file will be visible in the Testing tab.

In the Test the imported failing test file should FAIL, that proves our rule is precise (filtering irrelevant documents) and working correctly.
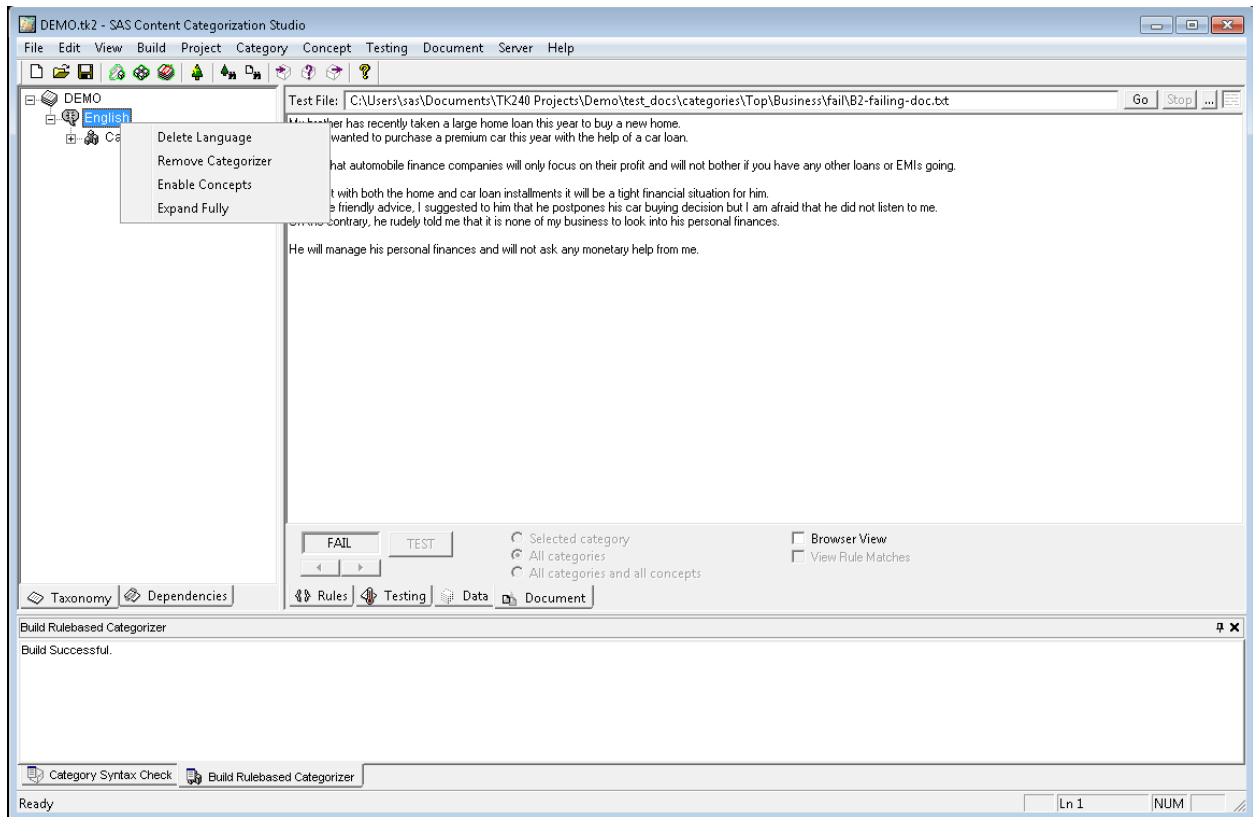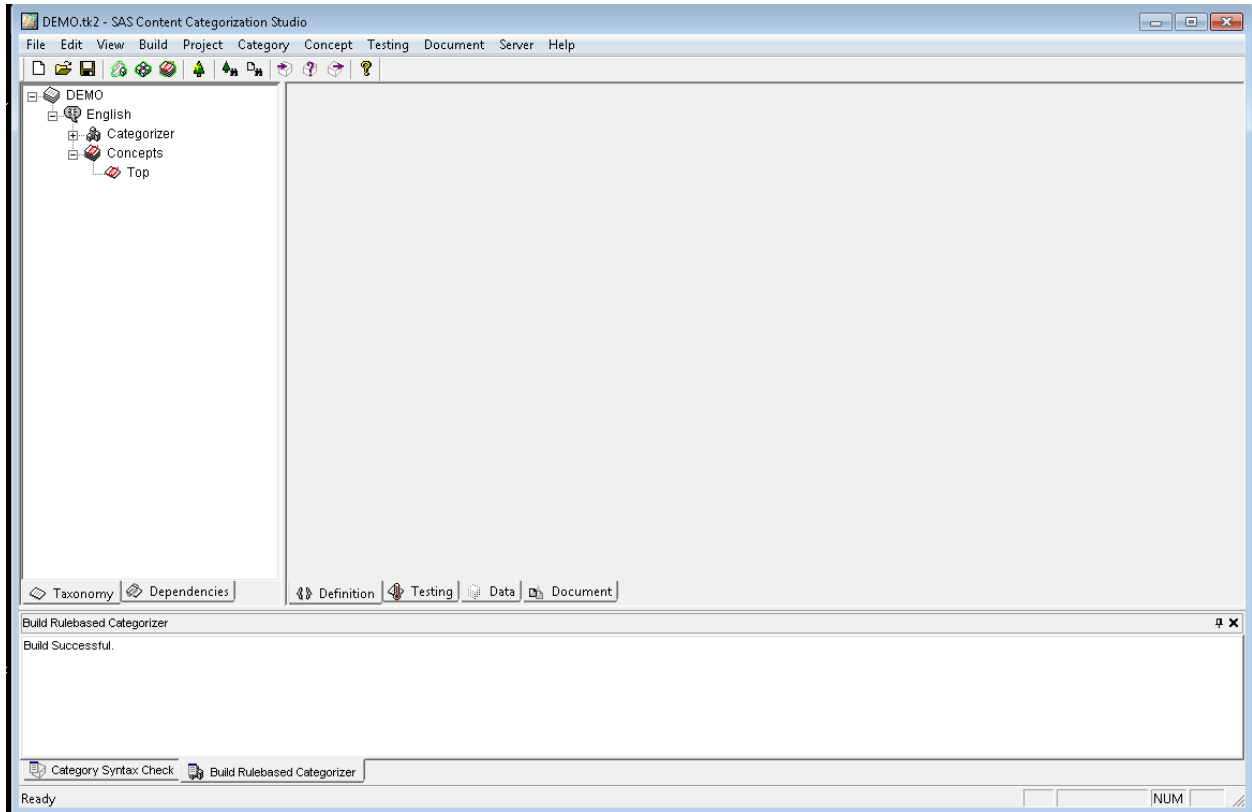


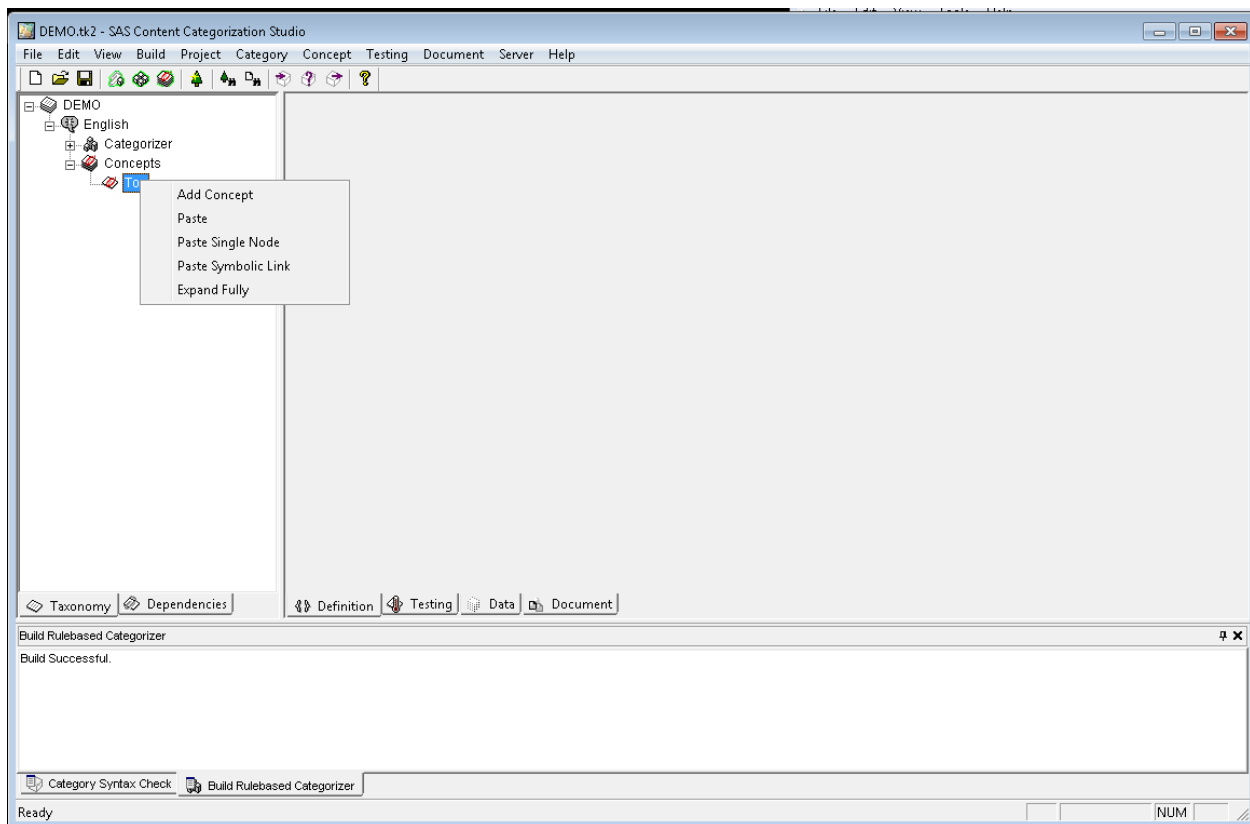11. Double clicking the Failing document opens it in the Document Tab.

## Part 3:  Creating LITI Definition Concepts

1.  Right mouse click on English and then click Enable Concepts to create Concept Taxonomy
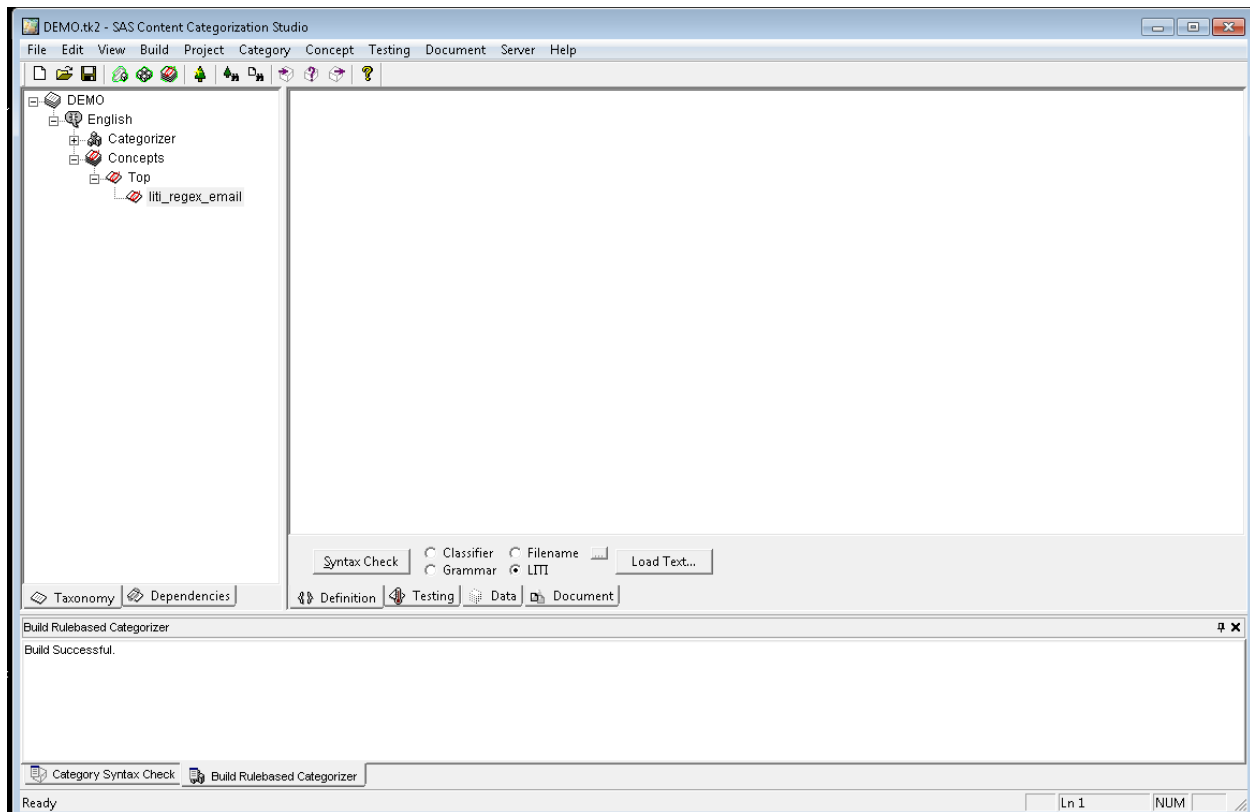
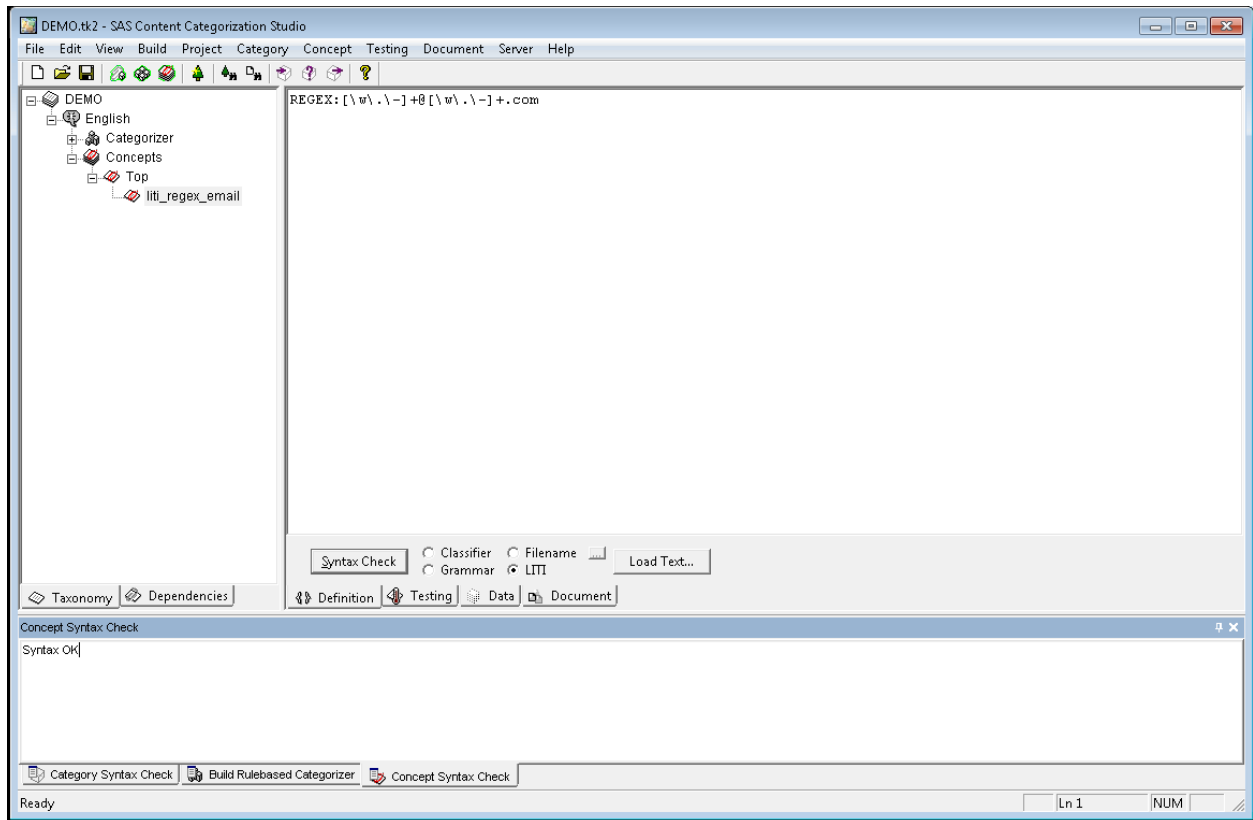2. To add new Concept do Right Click on Top and select Add Concept

3. Create a concept definition named liti_regex_email. In the below screenshot the LITI concept uses Regular Expressions to capture the email address ending with .com.  To create the LITI concept, select the LITI radio button on the Definition tab.
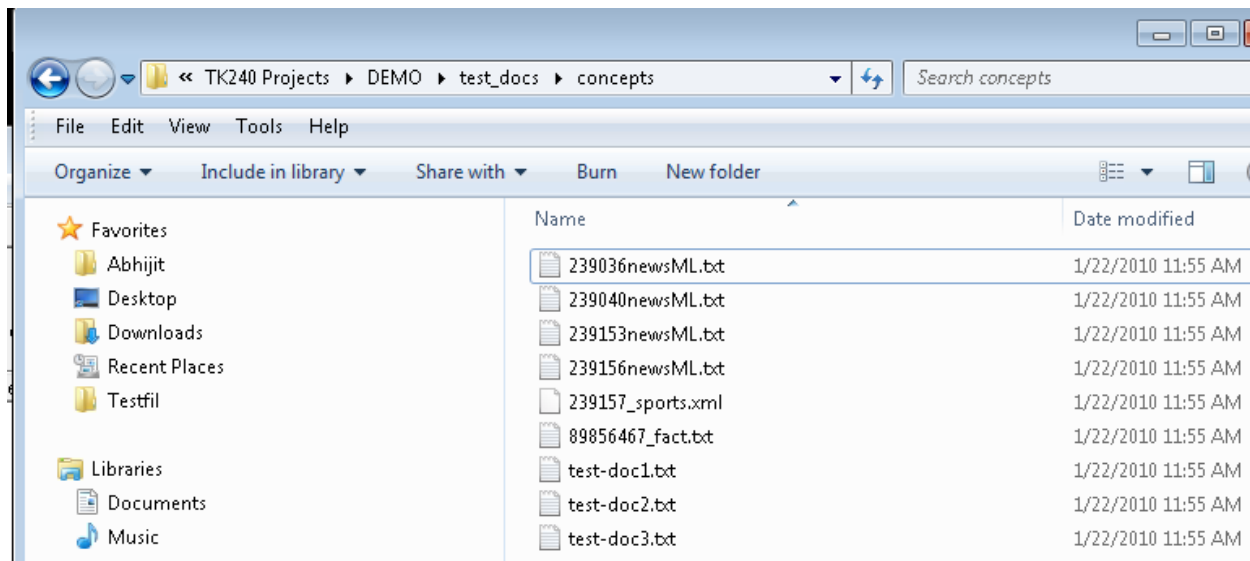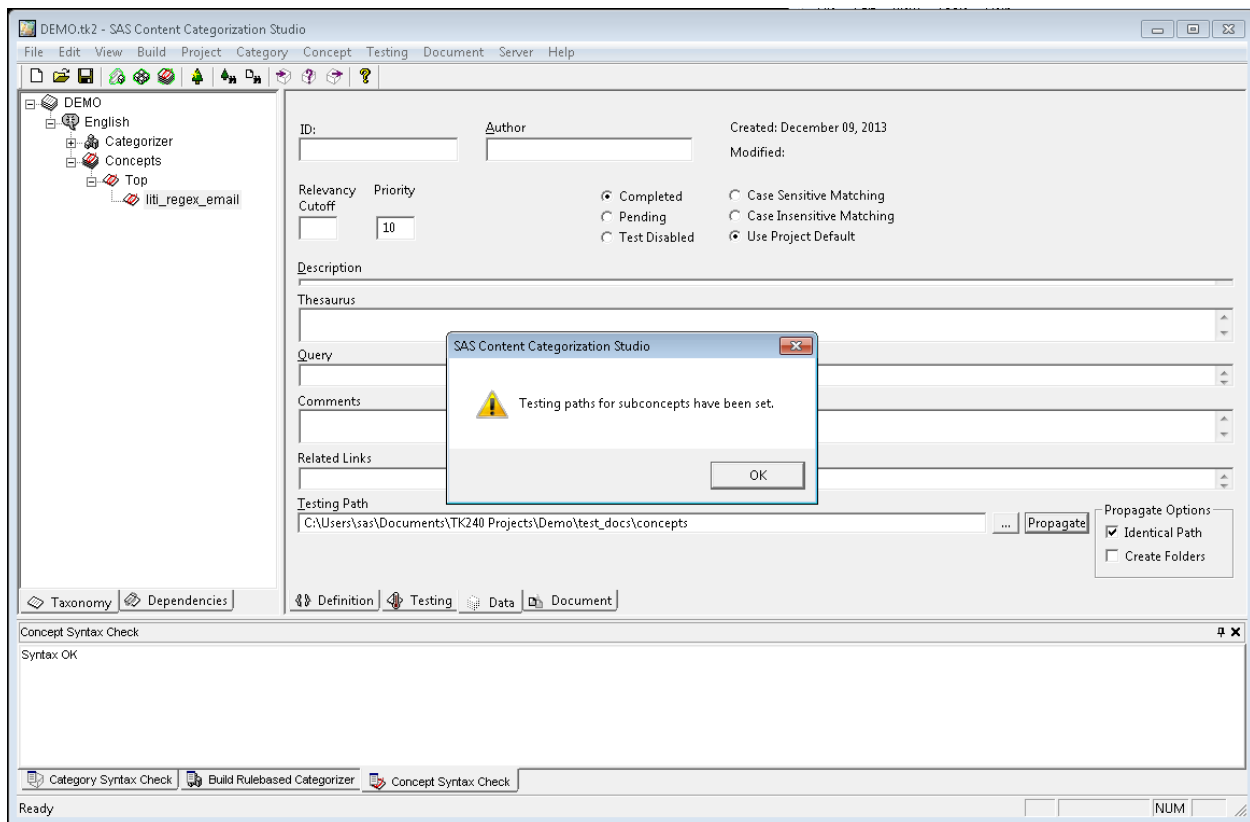
Write the definition
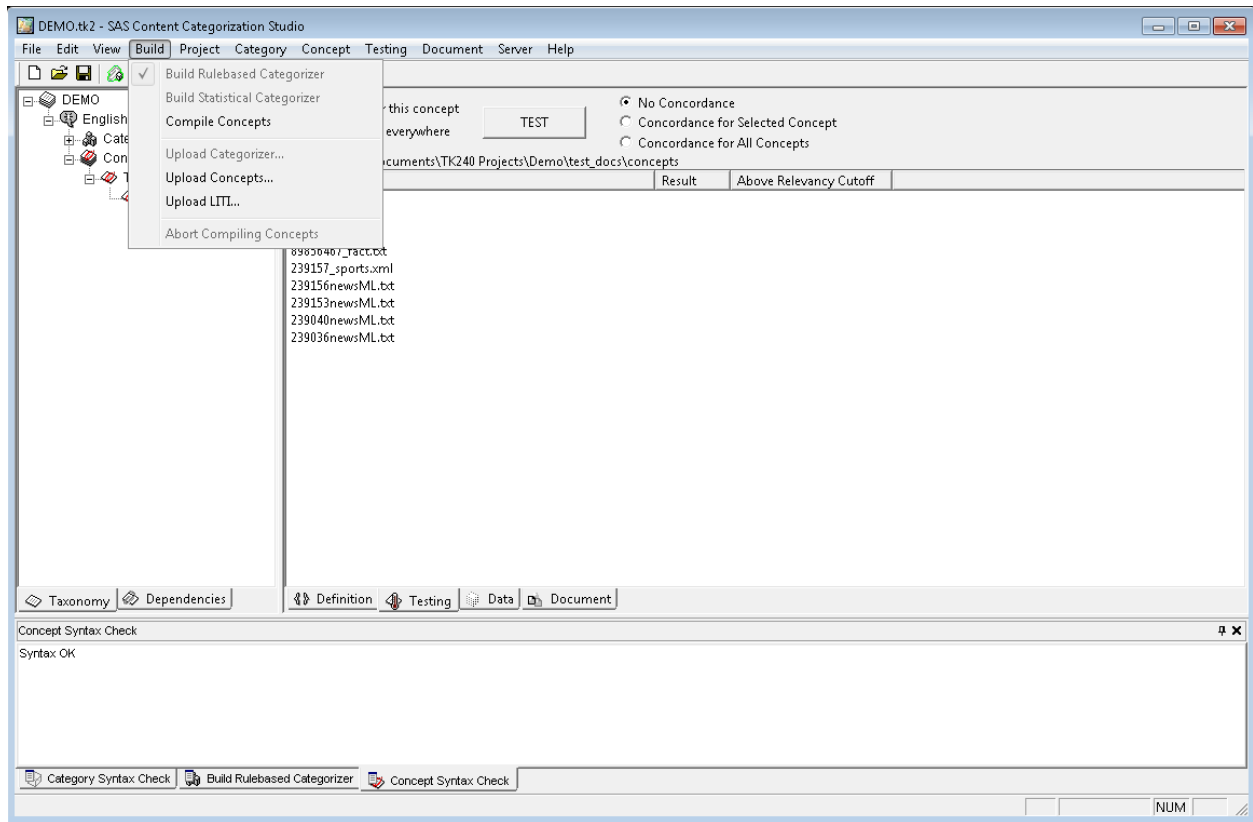
```
REGEX:[\w\.\-]+@[\w\.\-]+.com
```

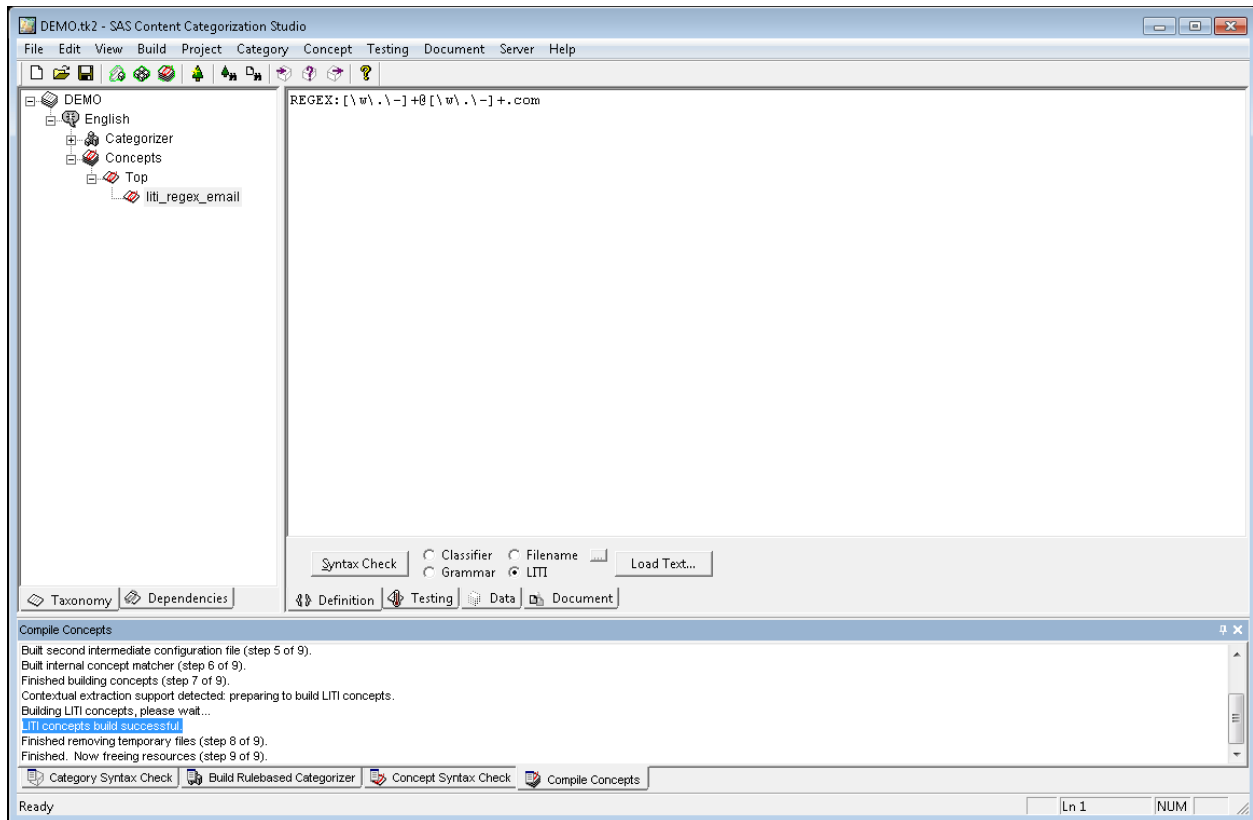on the Definition tab and then check for the syntax by clicking Syntax Check button.

4. On the Data tab set the path of the folder where the files to be tested are located and then click Propagate.
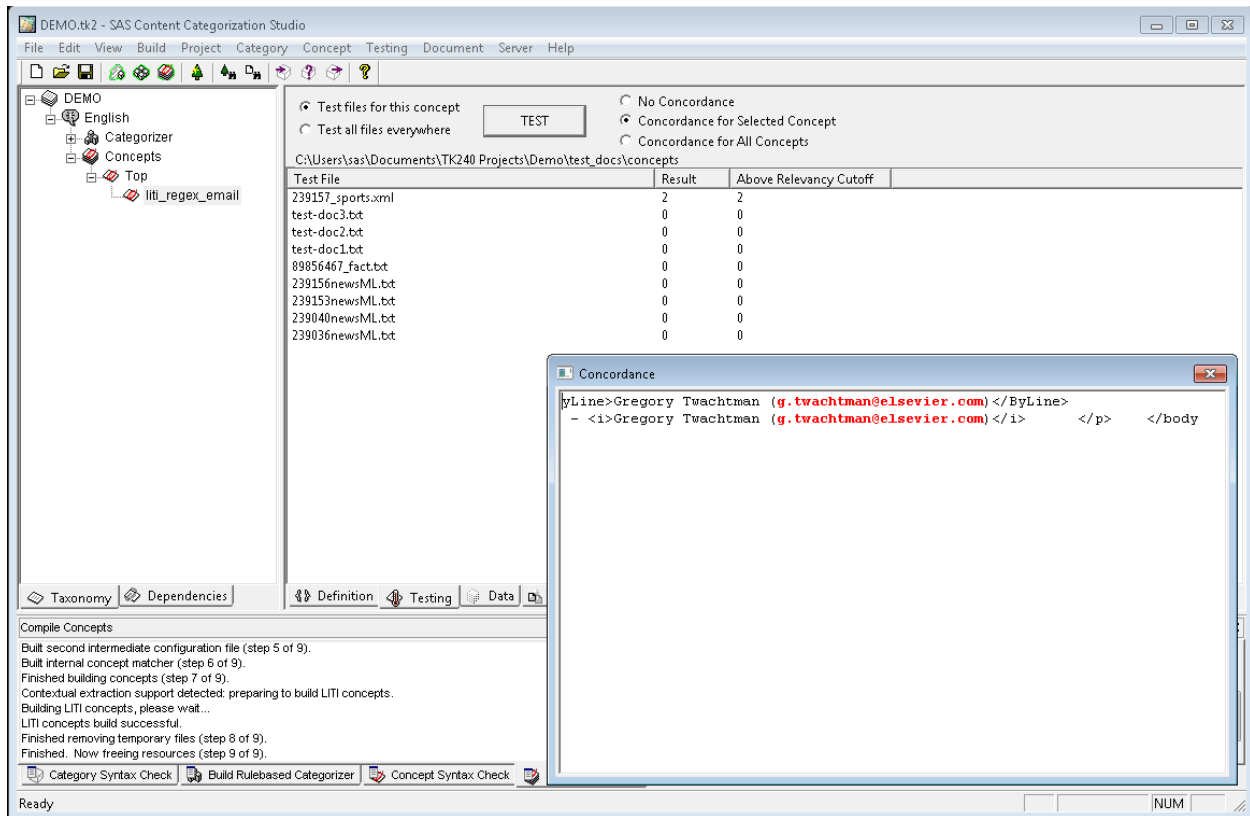
Files in the Testing Folder will now be visible in the Testing Tab.

5. Compile the concepts by clicking Build->Compile Concepts on the menu bar

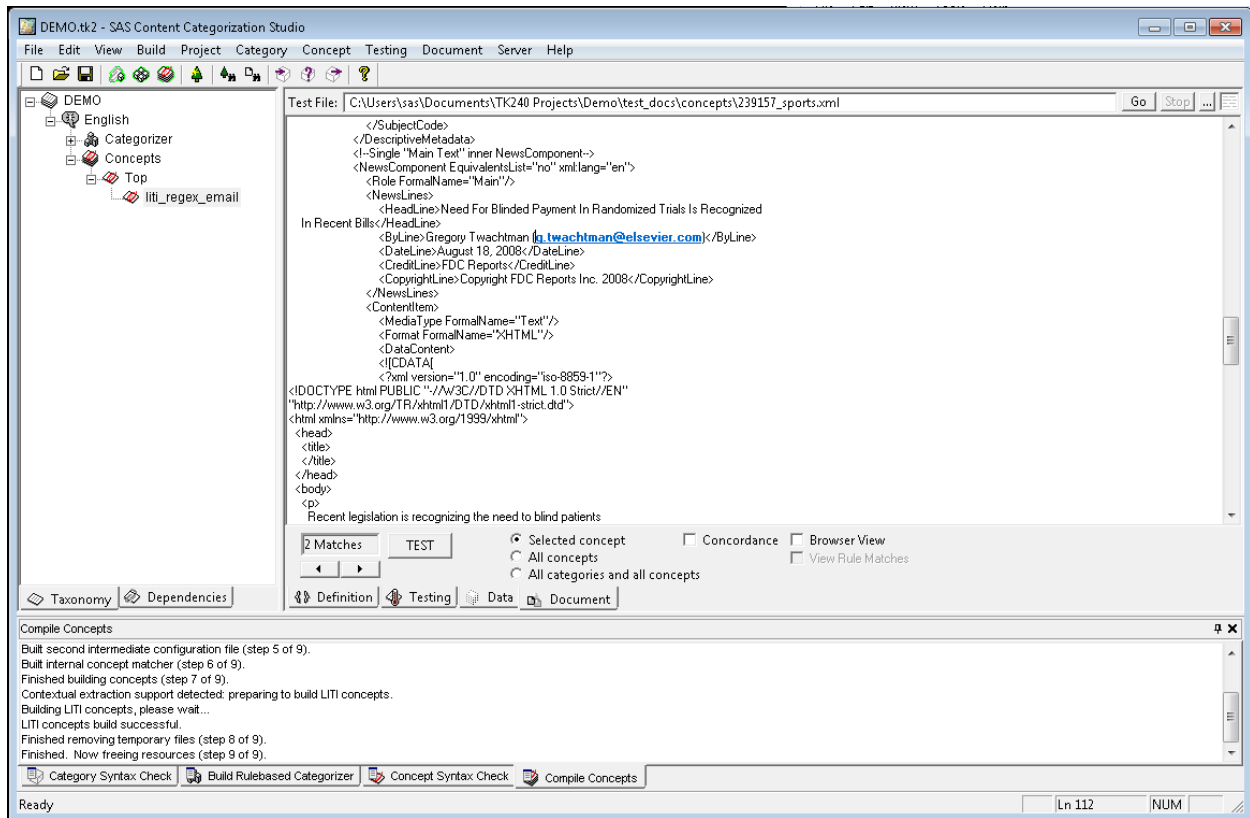LITI concept build is successful.

6. Now test the files for this concept by clicking the Test button on the Testing tab. With Concordance selected, we can control the number of characters, words or sentences displayed before and after each match.
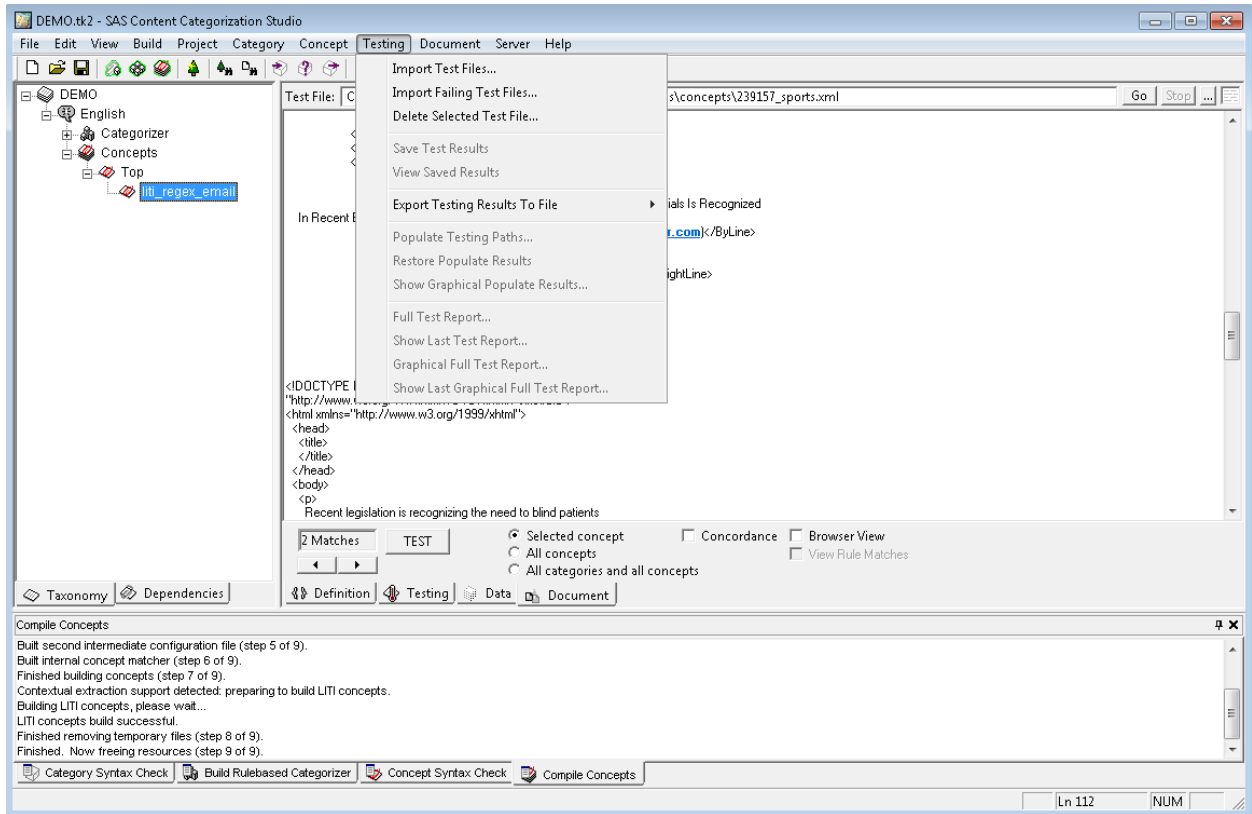
7. We can double click any of the test documents shown above to see which strings get matched. Matched strings will get highlighted in red as shown below.

Double click the 239157_sports.xml to open the file on Document tab and check which words get highlighted.
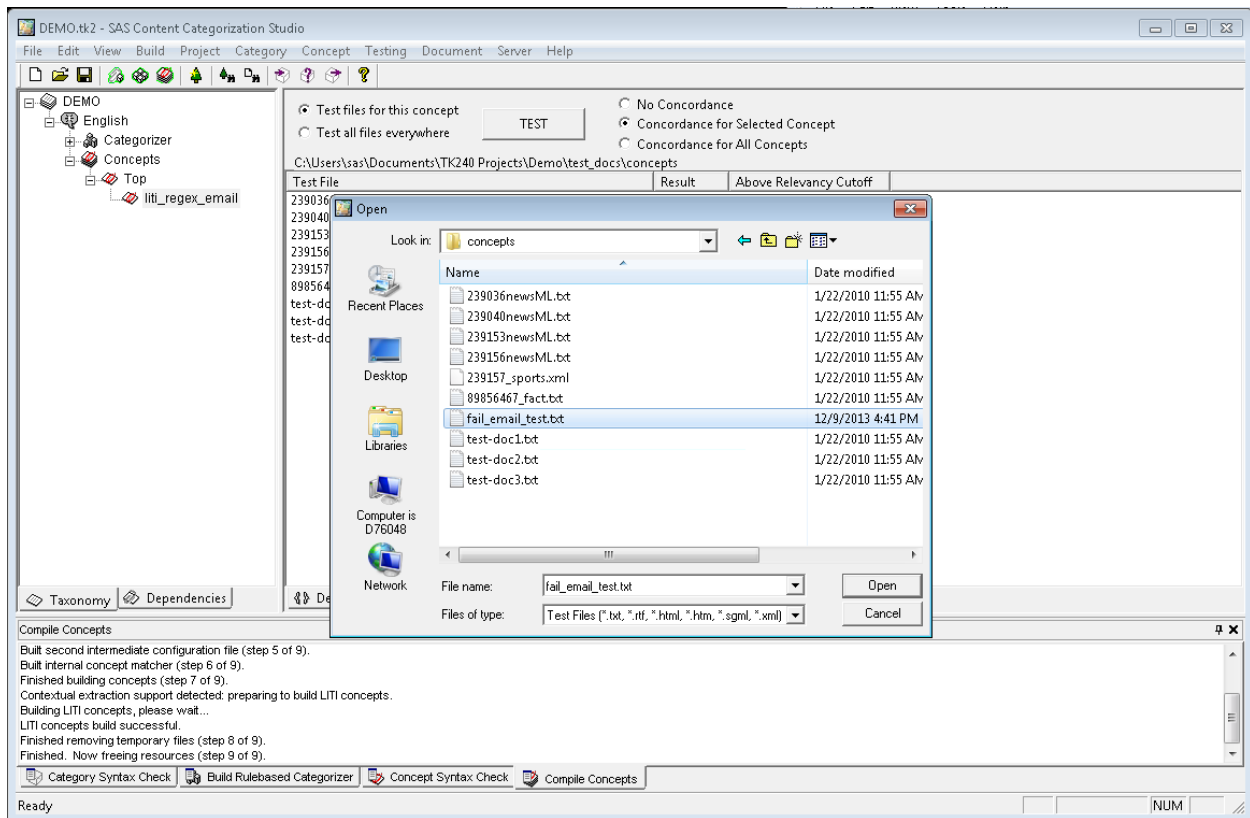
A good model should match all the relevant test documents (recall) while not matching irrelevant test documents (precision).
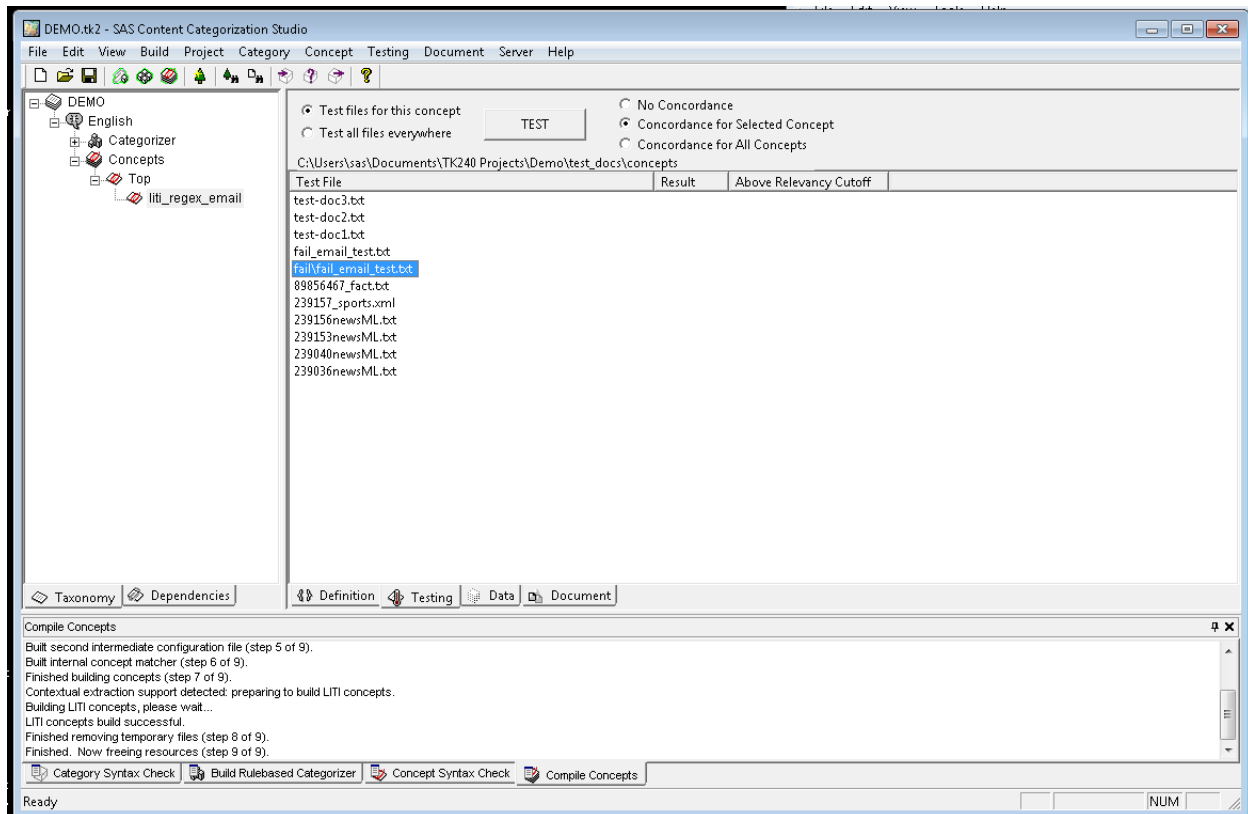
8. To test the failure of such irrelevant documents which do not satisfy the LITI definition of the concept, we can check by clicking Testing->Import Failing test Files… on the menu bar.
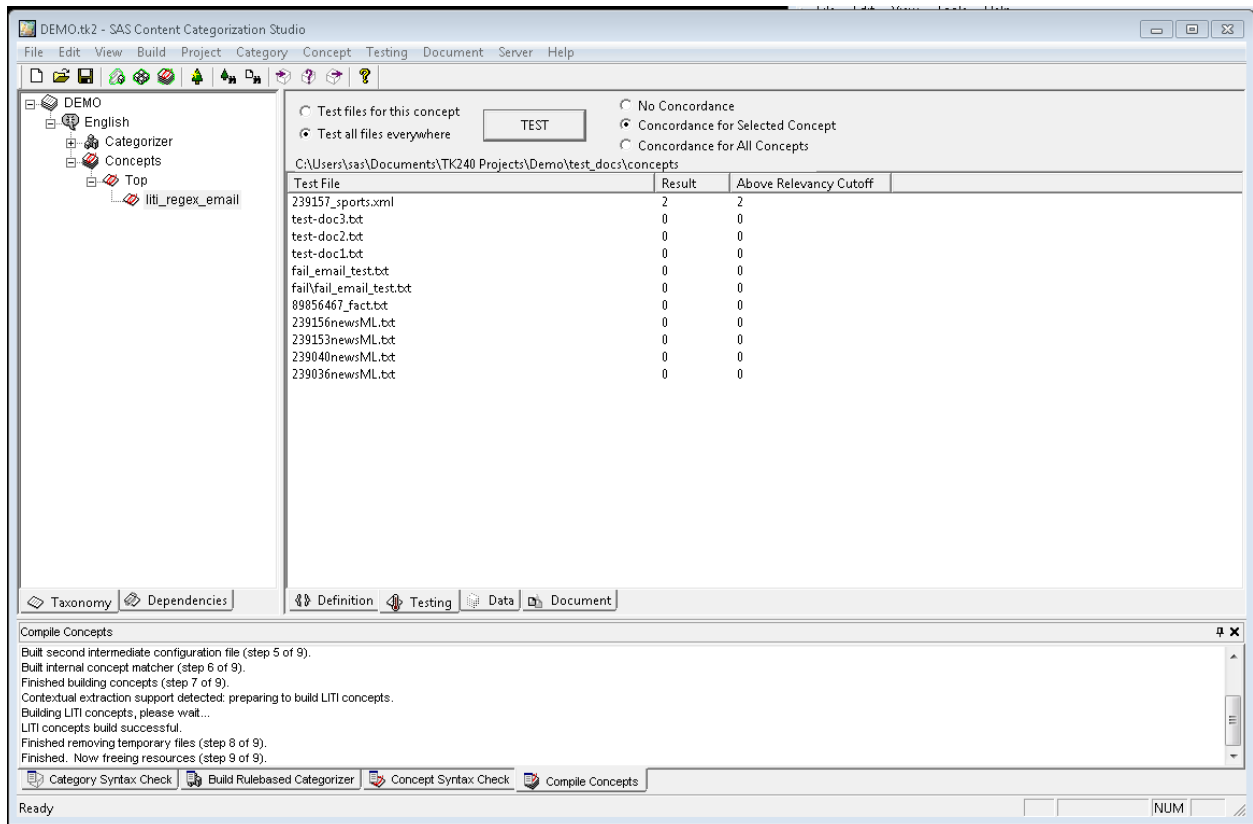
Open the failing document *fail_email_test.txt* to be tested as shown below

The imported file will be visible in the Testing tab.

In the test, the imported file fail_email_test.txt should not match any strings. You can verify this by noticing that the value for the Result column is 0. This shows that our concept definition is precise (filtering irrelevant documents) and working as expected.

No strings get highlighted for the file fail_email_test.txt as it does not satisfy our REGEX LITI Definition