![IBM logo]

developerWorks®

# Perform text analytics with Bluemix analytics services

**Lucas Silva**
Software Engineer
Big Data University

**Felipe Mosquetta**
Business Analyst Intern
Big Data University

**Marcos Rogério De Mello**
Software Engineer
Big Data University

25 February 2015
(First published 21 August 2014)

This tutorial explains how to perform text analytics using the Analytics for Hadoop service and dashDB service (formerly known as the Analytics Warehouse service) on IBM Bluemix™. Most of the processes are performed on a client machine with Eclipse IDE and the BigInsights™ plugin installed. After extracting the desired text, we use R in dashDB to plot charts with the results.

## Introduction

*Sentiment analysis,* or simply finding patterns in text, is one key tool many companies are using today to understand customers, monitor the reputations of their brands, or see what competitors are doing. Twitter and Facebook data are often used as input to perform such analyses. In this tutorial, we show you how you can easily gain valuable insight by performing such analyses. The 2014 Emmy Awards was used as the topic for our case scenario. Two lists with all the nominees were obtained from the USA Today web site. The first list includes the actors, and the second includes the series. These lists will be used to create our dictionaries to perform text analytics.

What you'll need for your application:

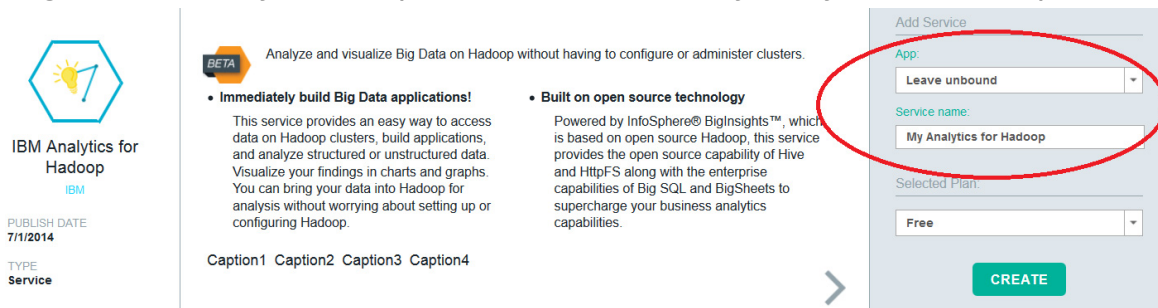- Eclipse Classic 4.2.2
- An IBM Bluemix™ account

## Setting up the Eclipse environment to perform text analytics

Assuming you already have Eclipse 4.2.2 installed on your computer, you need to add or install the InfoSphere® BigInsights Eclipse tools plugin:
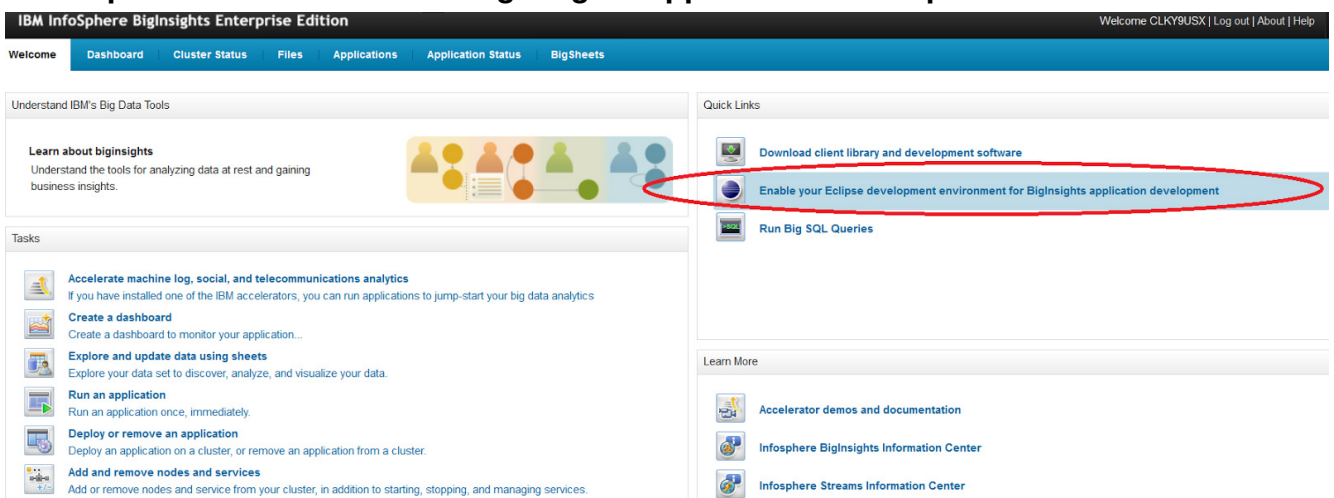
Trademarks

1. Log in to your Bluemix account. On the menu, click **Catalog** and scroll down until you find the Big Data section. Then choose **IBM Analytics for Hadoop**.
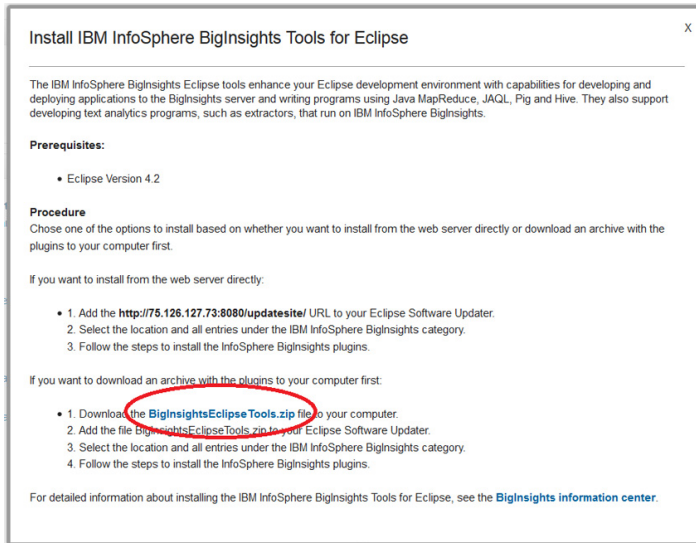


2. In the open window, click on the combo box **App** and select **Leave unbound**. Feel free to give a name of your own preference; we used "My Analytics for Hadoop." Click **Create**.
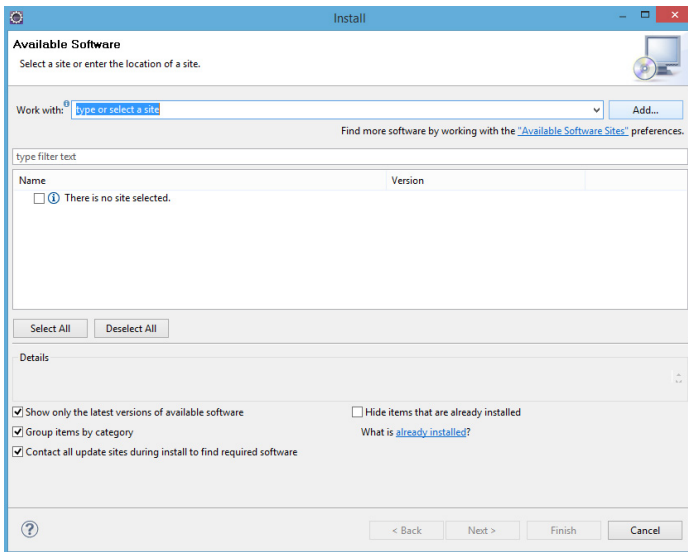


3. You will automatically be redirected to the service page. Click **Launch**.
4. You'll be sent to the InfoSphere BigInsights Enterprise Edition welcome page. In the Quick Links section, click **Enable your Eclipse development environment for BigInsights application development**.



5. You have two options to install the InfoSphere BigInsights Eclipse tools: install from the web server or from a .zip file. We will use the second one here so download BigInsigthsEclipseTools.zip.
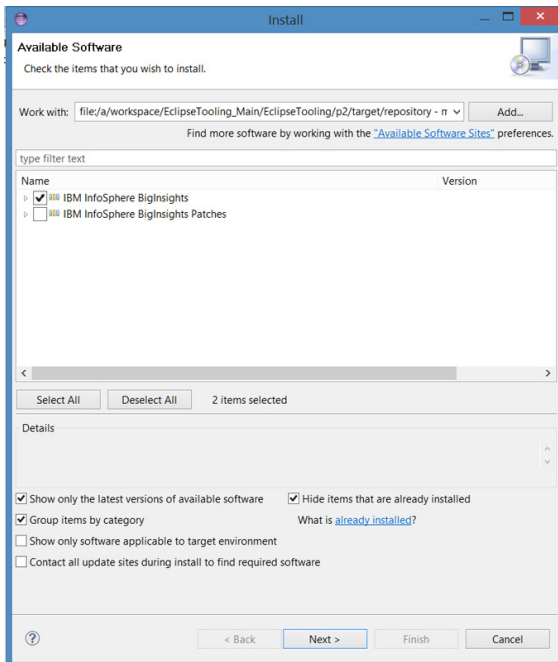
6. Now you need to install the plugin in the Eclipse IDE. Open Eclipse, select **Help > Install New Software**.



7. Click **Add** next to the **Work with** field, then select the BigInsightsEclipseTools.zip file (the file you downloaded) from your computer and click **OK**.

8. Check the InfoSphere BigInsights category and all of its features, then click **Next**. (Ensure that group items by category is checked so you see the category and its features).
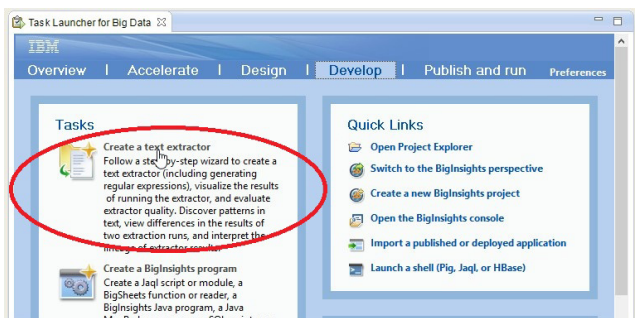


9. Follow the installation instructions.
10. Click **OK** if a security warning message appears.
11. Restart Eclipse when prompted.

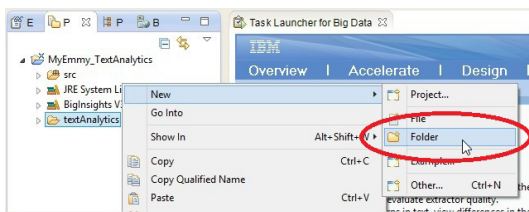# Creating a text analytics project

You need to create a new text analytics project on Eclipse IDE:

1. In the Task Launcher for Big Data, go to the Develop tab and choose **Create a text extractor** from the Tasks panel.
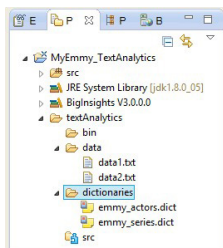


2. Name your project (suggestion: `MyEmmy_textAnalytics`) and click **Finish**.
3. Click **Yes** to switch to the InfoSphere BigInsights perspective when prompted.
4. Create a folder to put the data files. In the Project Explorer view, expand your project folder and find a folder called textAnalytics. Right click the textAnalytics folder

and select **New > Folder**. Name your folder (suggestion: `data`) and click **Finish**.



5. Now you can move your data to this folder. Select the folder you just created, copy your data files from your computer, and paste them in the Eclipse folder. You can also import the files with your data (**File > Import**) or drag and drop the file with your data from your local file system onto the data folder in Eclipse you just created.

6. Create a folder for the dictionaries. In the Project Explorer view, expand your project folder and select the **textAnalytics** folder. Right-click the textAnalytics folder and select **New > Folder**. Write `dictionaries` as the name of the folder (you can choose a different name, but this may result in divergences throughout this article) and click **Finish**.

7. Now you can move the dictionaries files to the folder that you just created. Two dictionaries contain the nominees of the International Emmy Awards: one with all the actors and actresses (emmy_actors.dict) and the other with all the series (emmy_series.dict). The next image shows the data, the dictionaries folders, and the files.
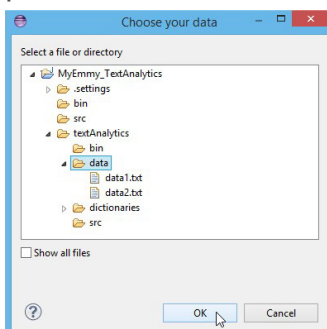


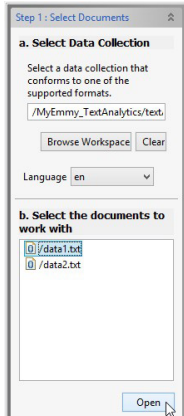Now you are ready to start the text analytics program.

## The text analytics workflow in Eclipse

In this section, we will go through the steps in the text analytics workflow in Eclipse to develop an extractor:
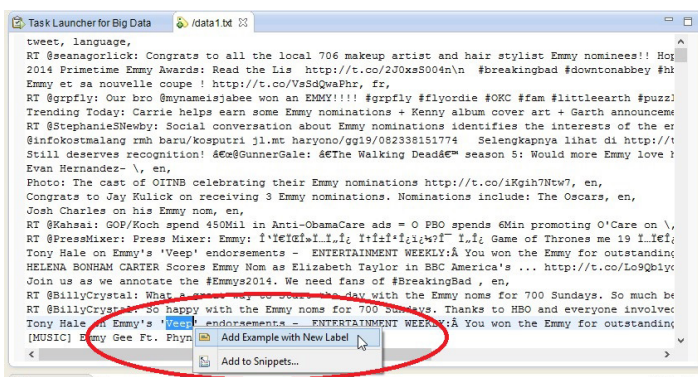
1. We need to specify the location of our data. In the Extraction Tasks view, go to **Step 1: Select Documents** and click **Browse Workspace**. Select the folder with your data (created in the previous section, called data) and click **OK**.
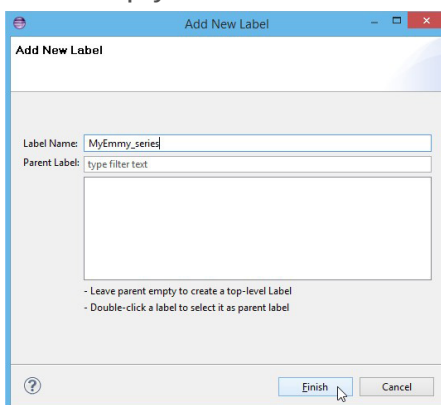
2. In the Language box, choose the language of your data. Select **en** for English. To use another language, you can check all the supported languages and their codes on the Supported languages webpage.

3. In Step 1b, select one of your data files and click **Open**. After this, you have completed Step 1 of the workflow.



4. In Step 2, you will label examples in the data. This step is useful in helping you to develop the extractor. In the data1.txt file, select **Veep**, right-click, then click **Add example with New Label**.
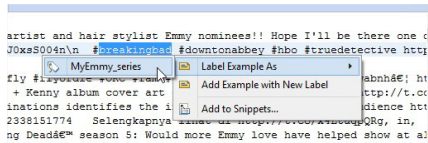


5. Name your label (suggestion: `MyEmmy_series`) in the Label Name field and leave the Parent label empty. Click **Finish**.
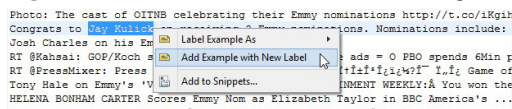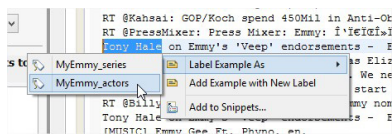
6. Select a different series (**Breaking Bad**, for example). Right-click, select **Label example as**, then click the name that you chose for your example in the previous step.
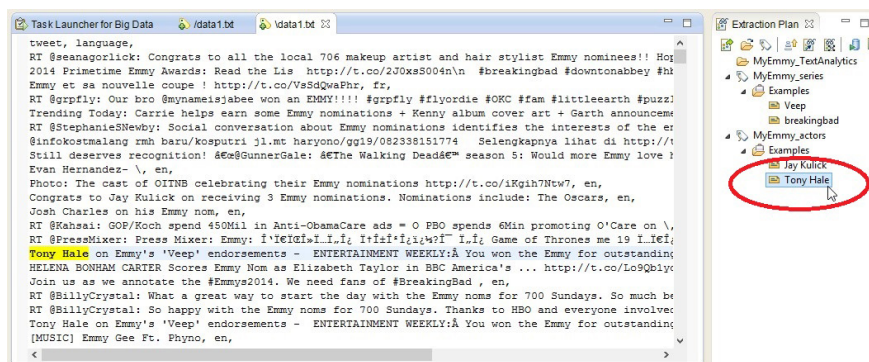


7. Let's create another label for the actors and actresses. Select **Jay Kulick** in the data1.txt file, right-click, then click **Add example with New Label**.
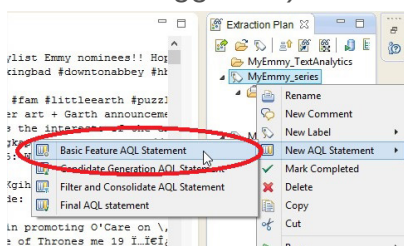


8. Write a name for your label (suggestion: `MyEmmy_actors`) in the Label Name field and leave the Parent label empty again. Click **Finish**.

9. Select a different actor's name (**Tony Hale**, for example), right-click, select **Label example as**, then click the name you chose for your example in the previous step.
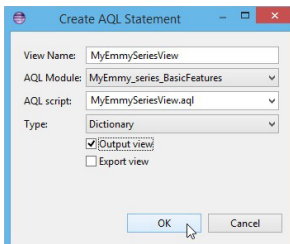


What we did in Step 2 was to select examples of actors and series present in the data1.txt file and assign them their specific labels. So we selected Veep and Breaking Bad as examples of series, and Jay Kulick and Tony Hale as examples of actors. We won't create more labels because we just need to find actors and series in our data. You can continue labeling the text, adding more examples of actors and series. At the end, look in the Extraction Plan to view the examples that you added. Double-click on one of them to highlight its location in the text.



10. In Step 3 of the text analytics workflow, you will create an extractor. In the Extraction Plan view, right-click on the name of your first label (**MyEmmy_series** if you gave the same name as suggested), select **New AQL Statement > Basic Feature AQL Statement**.
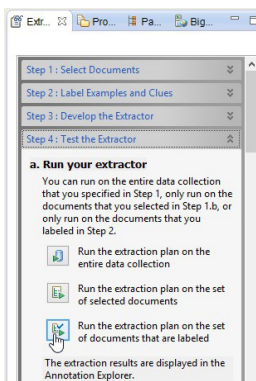
11. Give a name to your view (suggestion: `MyEmmy_seriesView`) in the View Name field and leave the AQL Module field with its default value. Then specify a name for your script (suggestion: `MyEmmy_seriesView.aql`) in the AQL script field. In the Type field, you have to select the type of your extract statement. To learn more, see The Extract Statement reference. Here, we will use dictionary as the type of our extractors, so we have to choose this option from the drop-down menu. Finally, select the **Output view** checkbox and click **OK**.



12. After that, an AQL file opens in the editor with some errors. We just have to fix them. We are going to start with the following code: `create dictionary MyEmmySeriesViewDict from file '<path to your dictionary here>' with language as 'en';`.

13. Replace `<path to your dictionary here>` for the path with the dictionaries with all the names of the series: `../../dictionaries/emmy_series.dict`. We need the `../` because the folder with the dictionaries is under the textAnalytics folder, and the current directory is the Emmy_series_BasicFeatures module. If you don't want to put it, just move the dictionary you want to use to your AQL file level. You can do this in the Project Explorer view. You can also create dictionaries from the examples labeled in the previous steps or add them to existing dictionaries. To learn more about this, read Step 10.b of Lesson 3: Writing and testing AQL.

14. Now we are going to work with this part of the code: `create view MyEmmySeriesView as extract dictionary MyEmmySeriesViewDict on R.<input column> as match from <input view> R;`.

15. Replace `<input column>` for text and `<input view>` for document and save your file. Then you should have code similar to the one below and with no errors.

```
module MyEmmy_series_BasicFeatures; create dictionary MyEmmySeriesViewDict
from file '../../dictionaries/emmy_series.dict' with language as 'en';
create view MyEmmySeriesView as extract dictionary
MyEmmySeriesViewDict on R.text as match from Document R; output view
MyEmmySeriesView;
```

16. We are going to test our extractor. There are three ways to run it: run the extraction plan on the entire data collection, on the set of selected documents, or on the set documents labeled. Select the third option: **Run the extraction plan on the set of documents that are labeled**.
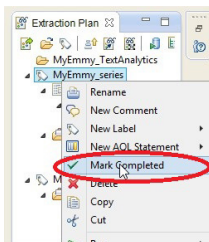
17. The results are shown in the Annotation Explorer view — specifically, in the Span Attribute Value column.
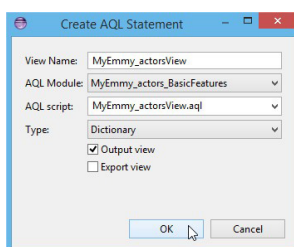


18. Now we can check this label as completed. Right-click on the name of the label and click **Mark Completed**.



19. Now, we need to repeat the same process to the MyEmmy_actors label. Do from Step 10 to Step 18 of this section for each label you have (in our case, we just have the MyEmmy_actors label left). The difference here is that you can create your view in the same AQL file of the previous one or create another aql file for your view. If you choose to include additional views in the same module, at the end, you will get just one tam file; otherwise, you will get one for each module. In this article, we will create the MyEmmy_actorsView in another module, which is the first option below.

- To create in another AQL file, in the Create AQL Statement window (obtained in Step 11), write a name for the view in the View Name field (suggestion: MyEmmy_actorsView), select MyEmmy_actors_BasicFeatures as the AQL Module, write a name for your AQL file, choose dictionary from the drop-down menu and, finally, check the Output view checkbox. Click **OK**.
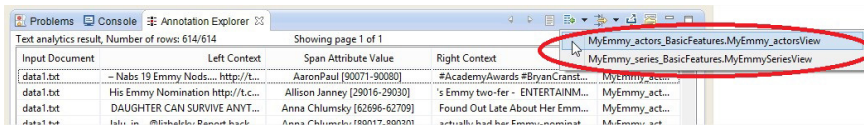


- To create in the same AQL file of the previous view, in the Create AQL Statement window (obtained in Step 11), select the same AQL module of the first label (MyEmmy_series_BasicFeatures) and the same AQL script (MyEmmy_seriesView.aql). Leave the other fields as shown in the preceding image.
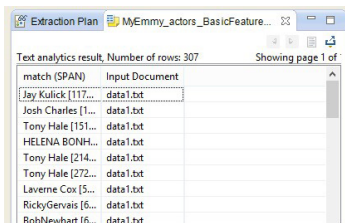
20. After creating your view and fixing the errors in the template (steps 12-15), you should have the following code.

```
module MyEmmy_actors_BasicFeatures; create dictionary MyEmmy_actorsViewDict
                from file '../../dictionaries/emmy_actors.dict' with language as 'en';
                create view MyEmmy_actorsView as extract dictionary
                MyEmmy_actorsViewDict on R.text as match from Document R; output view
                MyEmmy_actorsView;
```
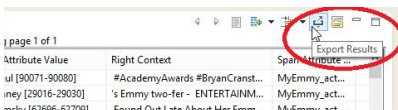
21. After running the extractor again (Step 16), you can select one of your views to see the results of each view separately. In this step, you can select the option to run your extraction plan on the entire data collection, unlike what was done in Step 16. Doing that, you will get results from all of your data files, not only from the labeled file.
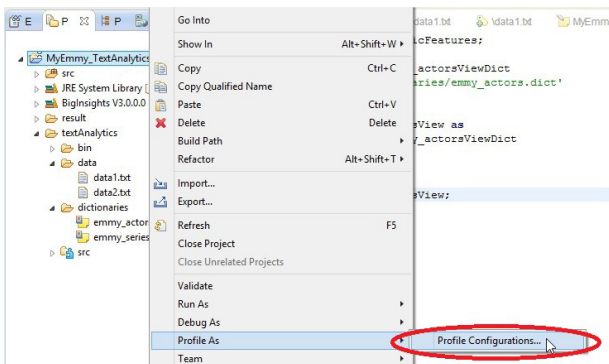


22. The results are shown in a different window beside the Extraction Plan view.



23. To finish Step 4 of the text analytics workflow, mark the MyEmmy_actors label as Completed (Step 18 of this section). Note: In this tutorial, we just use the Basic Feature AQL Statement, but you can also use the other ones. You can learn about them in Guidelines for writing AQL.

24. You can also export your data as HTML and CSV files after running the extraction plan on the entire data collection. In the Annotation Explorer view, click the **Export Results** icon.
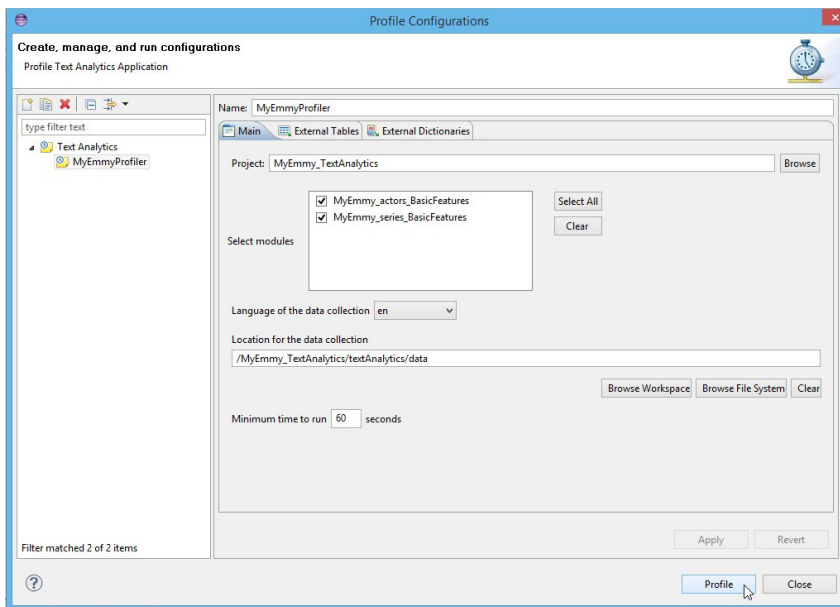


25. You now have to specify a path where the files will be saved. Click **Browse File System** and select the directory in which to save the results, then click **Finish**.

26. In Step 5 of the text analytics workflow, we can investigate the performance of our extractor. Go to the Project Explorer view, right-click on the name of your project, then go to Profile As and select **Profile Configurations**.
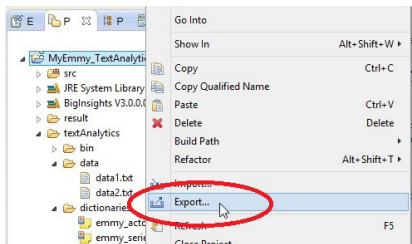


27. Double-click on **Text Analytics**, name your profile in the name field, select your modules, and specify the location for the data collection. You can also change the minimum time to the profiler runs. Then click **Profile**. After the profiler runs, you will get the performance runtime
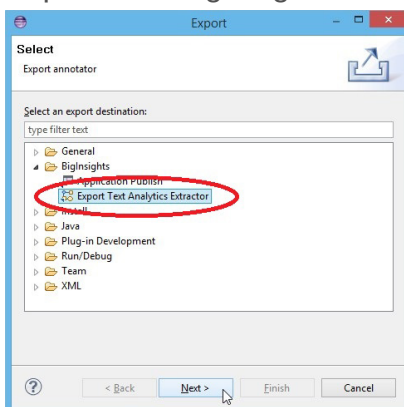
results of your extractor. To learn more about the AQL Profiler, see AQL Logic optimization.



28. The last step in the workflow is exporting the extractor. In the Project Explorer view, right-click on the name of your project, then go to Export.
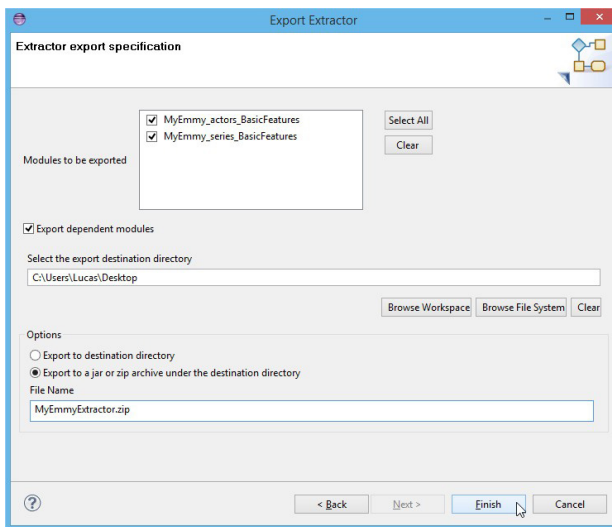


29. Expand the BigInsights folder, select **Export Text Analytics Extractor** and click **Next**.



30. Select the module to export and the **Export dependent modules** checkbox. Click **Browse File System** and specify the directory to save your extractor. In Options, select **Export to**

**a jar or zip archive** under the destination directory and write a name in the file name field.
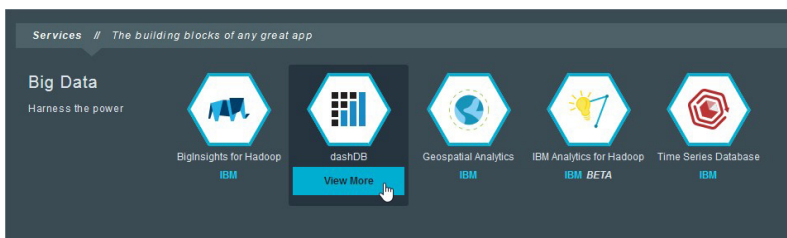


You should receive a confirmation message indicating that the process was successfully done. The steps in the text analytics workflow in Eclipse are complete. Additionally, you can upload your exported extractor to the HDFS on a Hadoop cluster to run it using much more data as input. To learn more, see the  BigInsights text analytics documentation .

## Analysis of results using R

We will use R included in the dashDB service in Bluemix to plot charts about the matches extracted in our data. Follow these steps to work with dashDB:

1. From the Bluemix dashboard, click **Catalog** in the menu, scroll down to the Big Data section, and click **dashDB**.

2. Change the name to `My dashDB` and click **CREATE**.



3. The dashDB service is created.
4. Click **LAUNCH**.
5. The dashDB service is running. As input, use the two CSV files obtained in The text analytics workflow in Eclipse. One file contains the list of actors and the other file contains the list of series. Load these datasets to dashDB. In **Load Data**, click **Load**.



6. Locate the input CSV files on your computer. The files do *not* include a header row. We show you how to do this for one file. Repeat the process for the other one.



7. Scroll down and click **Load File**.
8. Click **Next**. Select **Create a new table and load**.

9. Name the only column in the dataset as **MyColumn**, the table as **Actors**, and click **Finish**. The name of your table cannot have spaces. You will need this name soon.



10. With the data loaded to dashDB, let's perform some analysis using R. On the menu, choose **Analyze > Develop R Scripts**.



11. Now that you know how to load data, repeat the process to get the dataset with the series' name.

12. Upload the R Script. Click **Import** and locate the file in your computer.



13. The script is loaded and shown on the screen. Edit the script so the connection string uses the appropriate user ID and password. This information can be obtained from the Bluemix VCAP services. In the DashDB menu, click **Set Up > Connect Applications**. In the Connection settings section, the Bluemix VCAP services information includes the user ID. Copy this user ID.

14. Go back to the R script and paste the username in the appropriate location (on the fourth line in the following screen capture. (In this case, the user ID is dash101859). Between quotes, put the name of the table.

```
Save    Submit    Add a Data Frame...

library(ggplot2)
library(bluR)
mycon <- bluConnect("BLUDB", "", "")
Actors<- bludf(mycon, 'select * from BLU00068."Actors"')
Actors<- Actors$MyColumn
Actors<-as.factor(Actors)
table(Actors)
barplot(table(Actors),las=2,col=rainbow(46),cex.axis=1.0,cex.names=0.61,main="Tweets' Frequency of Actors
2014")

mycon <- bluConnect("BLUDB", "", "")
Series<- bludf(mycon, 'select * from BLU00068."Series"')
Series<- Series$MyColumn
Series<-as.factor(Series)
table(Series)
barplot(table(Series),las=2,col=rainbow(21),cex.axis=1.0,cex.names=0.61,main="Tweets' Frequency of Series
2014")
```

For the actors information, the fourth line of code reads:

```
Actors<- bludf(mycon, 'select * from dash101859, "Actors"')
```

For the series information, the eleventh line of code reads:

```
Series<- bludf(mycon, 'select * from dash101859, "Series"')
```

15. Click **Save** and then **Submit**. The Console Output shows a list of actors:

```
Actors
        "Aaron Paul "          "Allison Janney "           "Amy Poehler "
                 1                        1                          4
    "Andre Braugher "          "Anna Chlumsky " "Benedict Cumberbatch "
                 1                        3                          1
  "Billy Bob Thornton "          "Bob Newhart "         "Bryan Cranston "
                 2                        4                          3
    "Chiwetel Ejiofor " "Christina Hendricks "   "Christine Baranski "
                 1                        2                          1
        "Claire Danes "        "Fred Armisen "             "Gary Cole "
                 1                        2                          4
"Helena Bonham Carter "          "Idris Elba "             "Janefonda "
                 1                        2                          6
         "Jay Kulick "         "Jeff Daniels "            "Jim Parsons "
                 1                        4                          1
    "Joanne Froggatt "           "Joe Morton "              "Jon Hamm "
                 4                        4                          5
        "Jon Voight "          "Josh Charles " "Julia Louis-Dreyfus "
                 3                        3                         13
        "Julie Bowen "           "Kate Mara "          "Kate McKinnon "
                 1                        1                          2
       "Kate Mulgrew "   "Kerry Washington "          "Kevin Spacey "
                 3                       55                          4
       "Laverne Cox "        "Lizzy Caplan "           "Maggie Smith "
                28                        4                          2
     "Martin Freeman " "Matthew McConaughey "          "Mayim Bialik "
                 1                        3                          1
     "Natasha Lyonne "        "Paul Giamatti "         "Peter Dinklage "
                 1                        1                          4
      "Ricky Gervais "     "Taylor Schilling "             "Tony Hale "
                64                        2                          9
         "Uzo Aduba "
                41
```
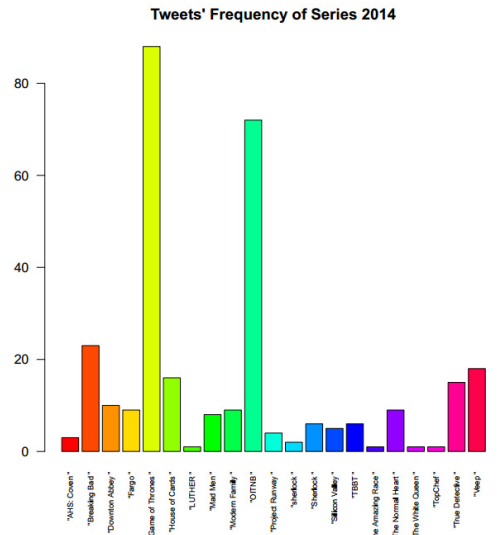
16. On the Plots tab is an icon for a PDF file that contains a plot of the desired information. Click the icon to download the PDF file. This file includes a plot of the Emmy series data and Emmy actors data. The following screen capture was cropped to show only the frequency of Tweets for the series nominations in 2014.

Our goal was to know whom people were cheering for in Emmy Awards 2014. So, according to this graph, we can easily observe three main crowds: *Game of Thrones*, *Orange Is the New Black*, and *Breaking Bad*. This graph shows that over 80 tweets were about *Game of Thrones*. We see, still by this graph, that people have expectations of the awards. For example, they do not care much about reality TV shows, such as *Top Chef* and *The Amazing Race*. It is good to remind you that this data was collected during the first four days after the official nomination.

# Conclusion

In this tutorial, we demonstrated how to create a text analytics program and some analytics based on a simple example of the Emmy Awards. You can use what you have learned here and apply in more complex examples based on your needs.

# Downloads

| Description | Name | Size |
| --- | --- | --- |
| Code file | codigo.txt | 1KB |
| CSV files | csvfiles.zip | 2KB |
| Data files | dataFiles.zip | 63KB |
| Dictionaries | dictionaries.zip | 2KB |

# Resources

## Learn

- Learn more about text analytics in the Knowledge Center.
- Learn more about the InfoSphere BigInsights text analytics.
- Learn more about the supported languages by annotators and dictionaries.
- Learn more about the extract statement.
- Learn more about writing and testing AQL.
- Learn more about the guidelines for writing AQL.
- Learn more about the AQL Profiler.
- Get Eclipse Classic 4.2.2 .
- Visit the developerWorks Information Management zone to find more resources for DB2 developers and administrators.
- Stay current with developerWorks technical events and webcasts focused on a variety of IBM products and IT industry topics.
- Follow developerWorks on Twitter.

## Get products and technologies

- Evaluate software in the way that suits you best: Download a product trial, try a product online, or use a product in a cloud environment.

## Discuss

- Get involved in the developerWorks community to connect with other developerWorks users while exploring the developer-driven blogs, forums, groups, and wikis.

# About the authors

### Lucas Silva

Lucas Lopes da Silva is a software engineer at Big Data University and fourth-year computer science student currently on a scholarship program at the University of Toronto.

---

### Felipe Mosquetta

Felipe Henrique Mosquetta Oliveira is a business analyst intern at BDU and fourth-year statistics student on scholarship at the University of Toronto.

---

### Marcos Rogério De Mello

Marcos Rogério de Mello is a third-year computer science student, an international student on scholarship by the Science Without Borders Program at the University of Toronto, and an intern at Big Data University.