**BigInsights Cloud Tutorial: Analytics for Hadoop on Bluemix**

**Tutorial 3: Loading data into BigInsights**
Learn how to quickly load data into your InfoSphere® BigInsights™ Hadoop environment, the options available to you and how to quickly review the data at hand.

Business data is stored in various formats and sources. Before you import your data into the InfoSphere BigInsights distributed file system, you must determine what questions you want to answer through analysis, identify the data type of your sources, and use the tools and procedures that best fit your business need.

You can use InfoSphere BigInsights with your existing infrastructure or data warehouse to import data and content in its original formats, or you can import huge volumes of at-rest (static) data or incoming data in motion (continually updated data). After you import your data, you can explore the data separately or combine the data to complete exploration and analysis.

Many businesses might want to examine the popularity of a specific brand or service in social media. The data that is provided for this lesson is the result of a BoardReader application search for the instances of the phrase "IBM Watson™" on the Internet. This search is detailed in the developerWorks® article, Analyzing social media and structured data with InfoSphere BigInsights: Get a quick start with BigSheets.

You are going to need this data as referenced in the above article, so please download it before proceeding. (The download is at the bottom of the article, yet feel free to [download directly, here](#))



Accept the terms and conditions and save the file article_sampleData to your local system. After you unzip the file, the article_sampleData folder should contain the files RDBMS_data.csv, blogs-data.txt, news-data.txt, and a README.txt file that details the data output.
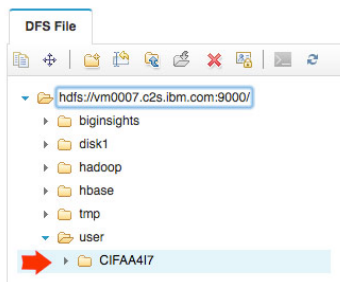
Make sure to note where you save these files – we will be uploading them to the cloud shortly.

For this tutorial, and the related tutorial on BigSheets (Tutorial 4 for our Bluemix series), only the news and blog data that was returned by the search is used. The returned data was slightly modified to contain only a subset of the information that the BoardReader application collects from blogs and news feeds. The full-text/HTML content of posts, news items, and certain metadata, was removed to keep the size of each file manageable. For information on BoardReader, [click here](#).

Now, let's look at the file system and create a folder where we will load the data.

**Creating a Folder**

1. Open the InfoSphere BigInsights Web Console.

2. From the **Files** tab, select the **HDFS** folder. (Note – the HDFS folder may already be open) – Within the **HDFS** folder navigate to the **USERS** folder and find your specific folder with your user id (note: this is the ID that is located on the IBM Analytics for Hadoop startup page). Navigate to this folder.



3. Create a directory to store this data in the distributed file system. Click the **Create**

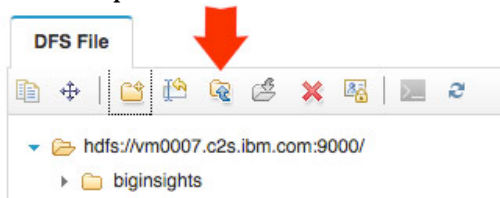**Directory** folder icon (  ) in the Files toolbar.

4. Name your directory. For this lesson, in the **HDFS** directory, create the directory bi_sample_data. For example the home directory is /user/<USERID>/.

You now have a directory to store all of your source data files and application results. Next, we will load the data.
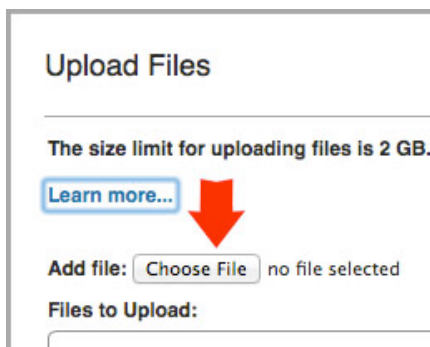
**Loading the data:**

Note: There are many ways to load data into BigInsights, but for this tutorial we are going to aim to keep it simple. We will be using the native Upload funciton in BigInsights, which is only recommended for smaller data sets (under 2GB). For larger data sets, you can leverage the [Distributed File Copy application](#).
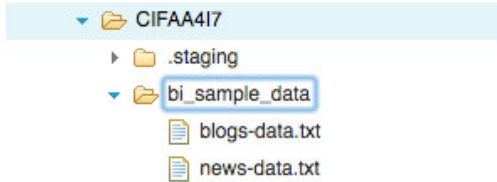
1. Make sure that you are in your 'bi_sample_directory' folder
2. On the files toolbar, click the **Upload..** icon – A new window should open 'Upload Files'



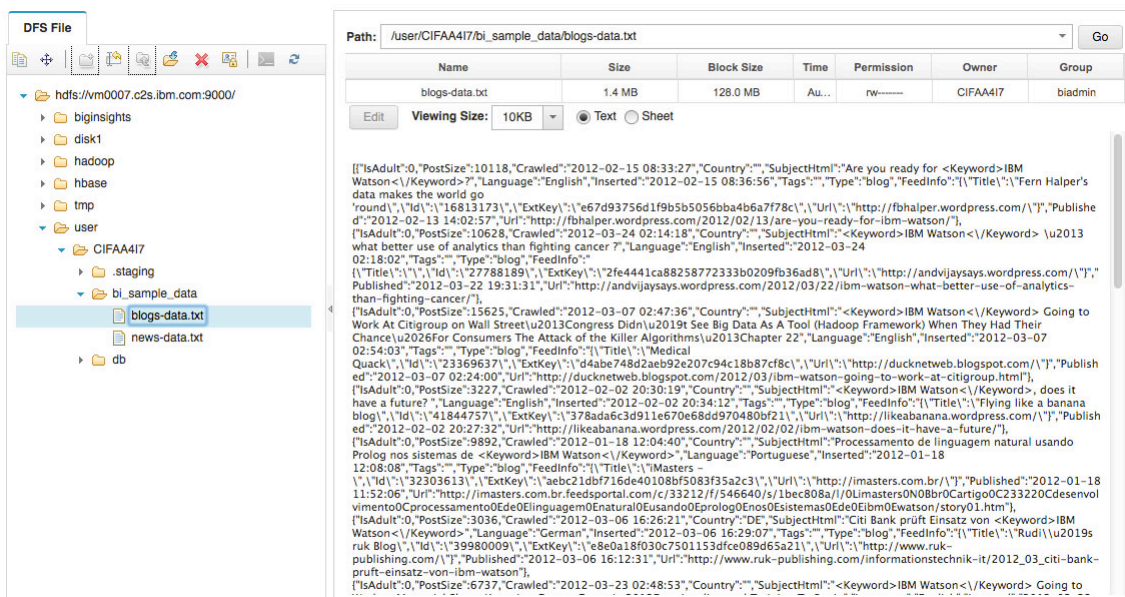3. Click on the **Choose File** button in this new window.



4. Navigate your local file system to where you saved the blogs-data.txt & news-data.txt files earlier in the tutorial.
5. Choose blogs-data.txt and upload
6. Repeat the process and Choose news-data.txt
7. Both files should be in the **Upload Files** window. Click the Ok button.
8. Navigate back to your 'bi_sample_directory' folder to ensure that the files are in there. Note: you can hit refresh on the toolbar to refresh the view.

You have now loaded the data into BigInsights that we will be using for the 4th tutorial.

If you want to take a sneak peak at what you just loaded – click on either of the files that we just brought in – the window to the right will show you a quick sample of the semi structured data.



We will tame and make sense of this data in the next tutorial, Exploring Data with BigSheets.