

Reading BBC news in Spark

1. Load the following file to Databricks

https://raw.githubusercontent.com/bbc/datalab-ml-training/master/datafest/bbc_news/tech/031.txt

2. Use Spark to read the file into an RDD
3. Use the `flatMap` function and the `split` function to split the content of the RDD into separate words and then count the total number of words in the file.
There should be 891 words if the text is split using `split()`, or 897 words using `split(' ')`.
4. Use the RDD `filter` method to count the number of times that the word **security** appears in the file (including the capitalised word **Security**). There should be 12 occurrences in total.