

Exercise Overview

In this exercise we will play with Spark [Datasets & Dataframes](#), some [Spark SQL](#), and build a couple of binary classifiaction models using [Spark ML](#) (with some [MLlib](#) too).

The set up and approach will not be too dissimilar to the standard type of approach you might do in [Sklearn](#). Spark has matured to the stage now where for 90% of what you need to do (when analysing tabular data) should be possible with Spark dataframes, SQL, and ML libraries. This is where this exercise is mainly trying to focus.

Feel free to adapt this exercise to play with other datasets readily availabe in the Databricks enviornment (they are listed in a cell below). #####Getting Started To get started you will need to create and attach a databricks spark cluster to this notebook. This notebook was developed on a cluster created with:

- Databricks Runtime Version 4.0 (includes Apache Spark 2.3.0, Scala 2.11)
- Python Version 3

Links & References

Some useful links and references of sources used in creating this exercise:

- Note:** Right click and open as new tab!
- [Latest Spark Docs](#)
 - [Databricks Homepage](#)
 - [Databricks Community Edition FAQ](#)
 - [Databricks Self Paced Training](#)
 - [Databricks Notebook Guide](#)
 - [Databricks Binary Classification Tutorial](#)

Get Data

Here we will pull in some sample data that is already pre-loaded onto all databricks clusters.

Feel free to adapt this notebook later to play around with a different dataset if you like (all available are listed in a cell below).

```
# display datasets already in databricks
display(dbutils.fs.ls("/databricks-datasets"))
```

Table						
	A ^B _C path	A ^B _C name	1 ² ₃ size	1 ² ₃ modificationTime		
41	dbfs:/databricks-datasets/sai-summit-2019-sf/	sai-summit-2019-sf/	0	1725169514125		
42	dbfs:/databricks-datasets/sample-jena/	sample-jena/	0	1725169514125		