Big Data Analytics - Lab 1

# Files, Shell, Bash, Scripting

On this first applied class we will start looking at files. The most basic unit to store information.

When dealing with very large amounts of data, visual tools like Excel or single host in-memory programming languages like python have a big performance degradation to the point of not working at all.

During the course of these lectures, we will learn how to work in environments capable of horizontally scale to datasets long above the 1M max rows of excel.

One of the original ways to use computers to process data was developed for the Unix Environment and followed its philosophy. Shortly summarized as:

- Write programs that do one thing and do it well.
- Write programs to work together.
- Write programs to handle text streams because that is a universal interface.
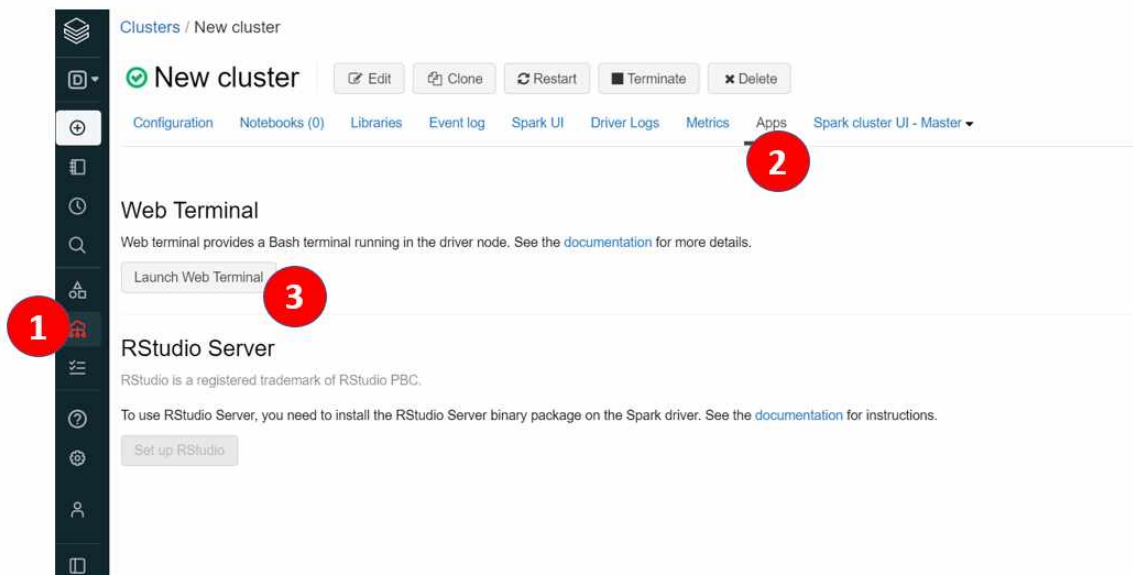
**Why starts this way?**
- Ease of execution of commands (no need to copy and paste every time)
- Powerful programming constructs

This mantra evolved extremely well, to the point that most of the commands developed within the Unix ecosystem, more than 30 years ago, are still relevant today. Not only that, but state of the art tools like Spark, which we will learn later in the course, follow a similar approach.

Today the course will be a brief introduction to Unix commands, often called Shell commands, and how they can be useful on modern day-to-day data processes.
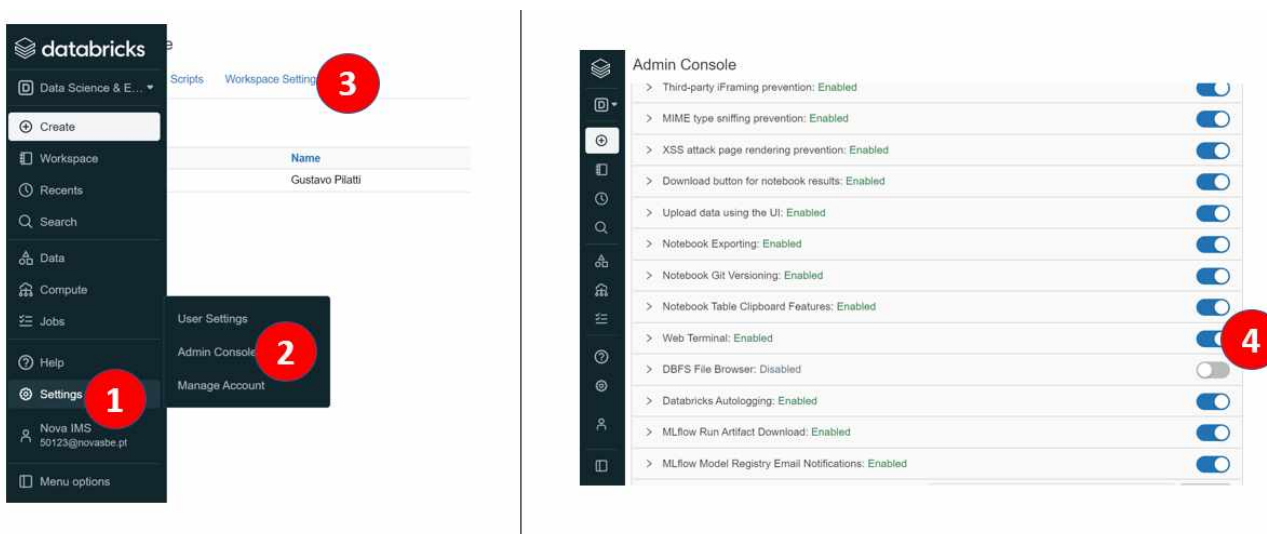
# Big Data Analytics - Lab 1

For the ones running MacOS or any Linux distribution you can run the class on your own computer. For windows users, Databricks runs on Linux, therefore all notebook exercises will work on Databricks.



If the Web Terminal is not enabled, follow these steps:

1. Go to the Settings
2. Click on Admin Console.
3. Click the Workspace Settings tab.
4. In the Advanced section, click the Web Terminal toggle.
5. Refresh the page.

Note: For the ones running windows it is also possible to run bash in it. Although the process is a bit more complicated but still possible. (https://docs.microsoft.com/en-us/windows/wsl/install-win10)

# Shell commands:

A Shell provides you with an interface to the operating system. It gathers input from you and executes programs based on that input. When a program finishes executing, it displays that program's output.

Shell is an environment in which we can run our commands, programs, and shell scripts. There are different flavors of a shell, just as there are different flavors of operating systems. Each flavor of shell has its own set of recognized commands and functions.

In this course we will use **Bash** ("Bourne Again Shell") shell as the main shell interpreter.

This shell is currently running on your local filesystem. Here is a list of shell commands to get comfortable with:

| | | |
|---|---|---|
| **ls** | ls [*options*] | Show directory contents, lists names of files. |
| **mkdir** | mkdir [*options*] directory | Creates a new directory of the specified name. |
| **cat** | cat [*filename*] | Display file's contents |
| **cd** | cd /*directorypath* | Change directory. Change to a certain directory name if provided. |
| **pwd** | pwd [-LP]<br><br>By default, `pwd' behaves as if `-L' were specified. | Displays the name of the working directory. |
| **touch** | touch *filename* | Creates a blank file with a specified name. |
| **less** | less [*options*] [*filename*] | View contents of specified file, page by page. |
| **head/tail** | tail [*options*] [*filename*] | Displays the first/ last 10 lines of a file. |
| **rm** | rm [*options*] *directory* | Removes a specified file. There is no recycle bin |
| **history** | history [*options*] | Display a listing of the last commands you've run. |
| **cp** | cp [*options*] *source destination* | Copy specified file to a new named file. Use -r flag copy a directory. |
| **mv** | mv [*options*] *source destination* | Rename a specified file or directory. |

# Big Data Analytics - Lab 1

| | | |
|---|---|---|
| **find** | find [-H] [-L] [-P] [-Olevel] [-D debugopts] [path…] [expression] | Search files and directories. Can use with wildcards (* ?[ ]). |
| **curl** | curl [*options*] url | Download a webpage |
| **help** | help [-dms] | Get help on a command eg. help ls |
| **echo** | echo [*options*] | Prints text to the terminal window |
| **grep** | grep [*options*] *pattern* [*filename*] | Used to search text for patterns specified by the user. |
| **uniq** | uniq [*options*] [*unput [output]*] | report or filter out repeated lines in a file |
| **sort** | Sort [*options*] [*filename*] | sort or merge records (lines) of text and binary files |
| **wc** | Wc [*options*] [*filename*] | Word count |
| **sed** | sed [options] | Pattern-matched string replacement |

**TIP:** You can press the up arrow to cycle through previous commands

**TIP:** When using windows, you can right-click to paste (instead of ctrl-v).

**TIP:**
- **[command] -h:** Display a file's help information.
- **[command] --help:** Display a file's help information.
- **whatis [command]:** Display a short blurb about the command.
- **Some commands do not have the --**help buit0in function, as echo, so you can enable it by typing: enable -n echo ; echo --help

**IMPORTANT:** CTRL-C (cmd-C) will cancel any command running, this will be useful if you accidentally try to open a large file.

# Let us start:

# Big Data Analytics - Lab 1

We start by understanding where we are within the file system and the content of the current directory.

hostname – what machine you are on

whoami – who you are logged in as

pwd - Show the current directory. Folders are divided with slashes "/".

ls - List the contents of the current directory.

Let us try to open a file with cat:

cat /etc/lsb-release - In this case we are opening a file called "lsb-release" in the "etc" folder that contains information about the operating system.

Now adding a bit of action, let us get a live stream of information.

htop - is an important command to check the status of the tasks running on the computer.





As a way to chain operations bash uses the concept operators. They are used to combine several operations together; this is where the power of

bash comes from. By chaining the simple commands shown in the previous section it allows for more advanced and useful operations.

## BASH operators:

| |
|---|
| **List terminators**<br><br>● "**;**": Will run one command after another has finished, irrespective of the outcome of the first. E.g.: `command1 ; command2` First `command1` is run, in the foreground, and once it has finished, `command2` will be run. |
| **Pipe operator**<br><br>● "**|**": The pipe operator, it passes the output of one command as input to another. A command built from the pipe operator is called a pipeline. E.g.: `command1 | command2` Any output printed by `command1` is passed as input to command2. |
| **Redirection operators**<br><br>These allow you to control the input and output of your commands. They can appear anywhere within a simple command or may follow a command. Redirections are processed in the order they appear, from left to right.<br><br>• "**<**": Gives input to a command. `command < file.txt` will execute `command` on the contents of "file.txt".<br><br>• "**>**": Directs the output of a command into a file. `command > out.txt` will save the output of `command` as "out.txt". If the file exists, its contents will be overwritten and if it does not exist it will be created.<br><br>• "**>>**": Does the same as "**>**", except that if the target file exists, the new data are appended. `command >> out.txt` If "out.txt" exists, the output of `command` will be appended to it, after whatever is already in it. If it does not exist it will be created. |
| **Multi-line execution**<br><br>Long commands can be separated into multiple lines by using a backlash after each line<br>`January February March April \`<br>`May June July August September October November \`<br>`December` |
| **Quotes** |

```
root@0206-124507-fis7e8yb-10-172-224-5:/databricks/driver# mkdir 'Documents and Settings'
root@0206-124507-fis7e8yb-10-172-224-5:/databricks/driver# ls
'Documents and Settings'   conf   eventlogs   ganglia   logs   metastore_db   preload_class.lst
root@0206-124507-fis7e8yb-10-172-224-5:/databricks/driver# mkdir Documents and Settings
root@0206-124507-fis7e8yb-10-172-224-5:/databricks/driver# ls
 Documents  'Documents and Settings'   Settings   and   conf   eventlogs   ganglia   logs   metastore_db   preload_class.lst
root@0206-124507-fis7e8yb-10-172-224-5:/databricks/driver#
```

**Tip**: attention to 'rm -r -f *'. For more information, follow this link. **Never run it***

# Example

Let's create a file and deal with it!

1. Create a file with the name "myfile.txt" with the words "banana apple carrot" as the file content

   echo "banana apple carrot" > myfile.txt

2. Verify if your file was created

   ls

3. All entries are in the same line. Let's delete this file

   rm myfile.txt

4. We need to create a file with one entry per line. Call this file 'data.txt'

   echo "banana" > data.txt

5. Append more lines with new entries

   echo "apple" >> data.txt | echo "carrot" >> data.txt |echo "watermelon" >> data.txt

6. Search for words with 'a'

   grep 'a' data.txt

7. Search for words with 'p' and 'c'

   grep '[pc]' data.txt

**Tip:** square parentheses are a matching set.

# Big Data Analytics - Lab 1
## High quality guide

http://www.compciv.org/bash-guide/