# LAB 1 INSTRUCTIONS

## DS8003 – MGT OF BIG DATA AND TOOLS

### Ryerson University

Instructor: Kanchana Padmanabhan

# Lab & Assignments

- Lab Computer
  - Username: same as ryerson
  - Password: same as ryerson
- Lab 1
  - Virtualbox
  - Hadoop setup
  - SSH login

# Lab 1 – Environment Setup

# Install VirtualBox

□ Download VirtualBox for your respective Operating System

   ▪ https://www.virtualbox.org/wiki/Downloads



   ▪ For MAC – file has extension ".dmg"
   ▪ For Windows – file has extensiion ".exe"
   ▪ Double click on the downloaded file and follow instructions

# Download Hortonworks HDP Sandbox

☐ Download HDP hadoop vm image: http://hortonworks.com/products/hortonworks-sandbox/

☐ The file will be called "HDP_2.3.2_virtualbox.ova" (or download the version available)

## Download & Install

The Hortonworks Sandbox provides an easy way to get started to learn and develop with the Hortonworks Data Platform (HDP) anywhere. You can either run it in the cloud or your personal machine.

### Hortonworks Sandbox on a VM

No data center, no cloud service and no internet connection needed! Full control of the environment. Easily extend with additional components or try the various Hortonworks technical previews. Always updated with latest edition.

**HDP 2.3.2 on Hortonworks Sandbox**
Runs on VirtualBox or VMware

Try out the very latest features and functionality in Hadoop and its' ecosystem of projects with HDP 2.3. Follow the Step by Step Tutorials.

System Requirements | Installation Steps | Release Notes
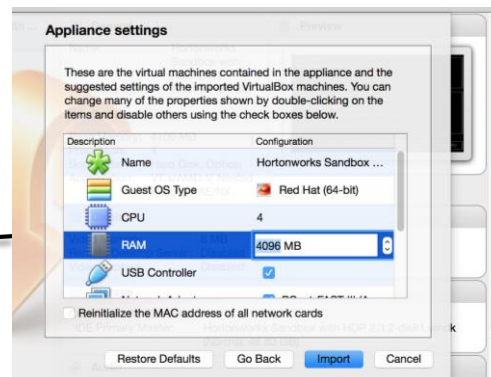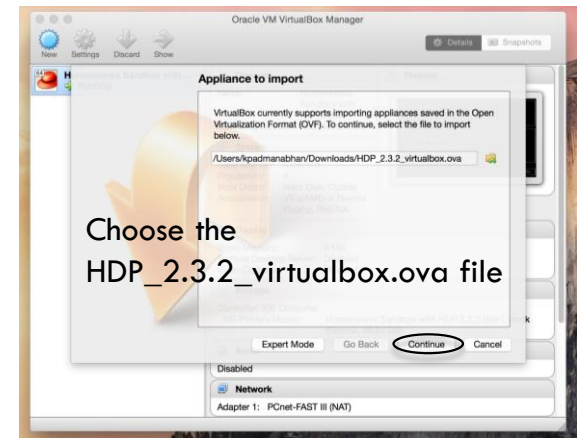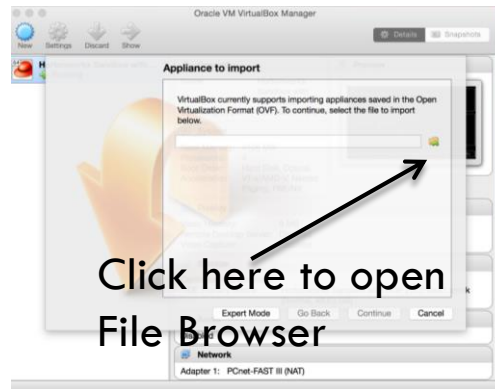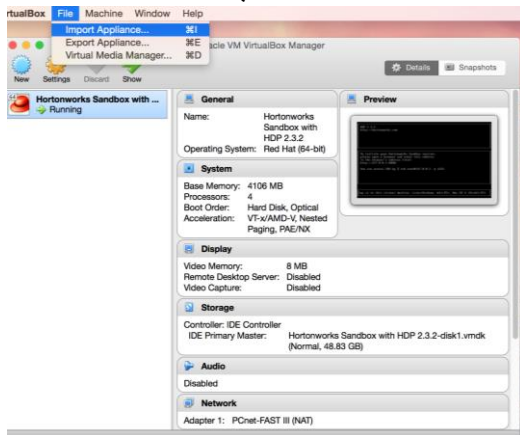
INSTALL GUIDES

for VirtualBox
Mac & Windows

for VMware
Mac & Windows

**for VirtualBox**
(HDP 2.3.2 - 8.5 GB)

**for VMware**
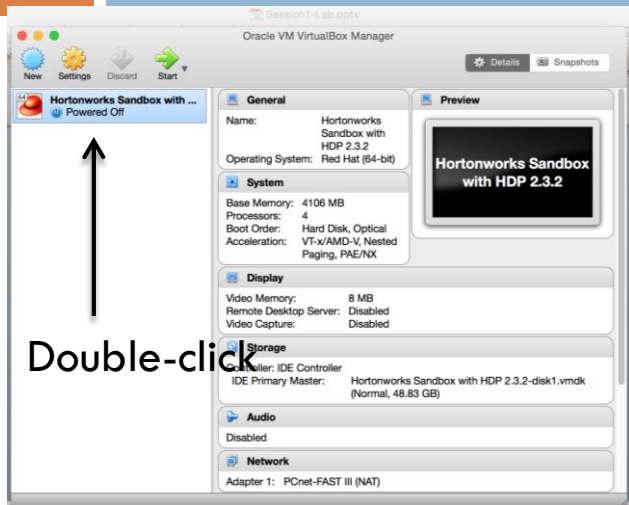(HDP 2.3.2 - 8.7 GB)

# Install Hortonworks HDP Sandbox

- Sandbox installation tutorial for Windows: http://hortonworks.com/wp-content/uploads/unversioned/pdfs/InstallingHortonworksSandbox2onWindowsusingVB.pdf

- Double-click on the virtual box

- Import the file "HDP_2.3.2_virtualbox.ova" (or latest version your downloaded) into the Virtual box (Follow pictures below)
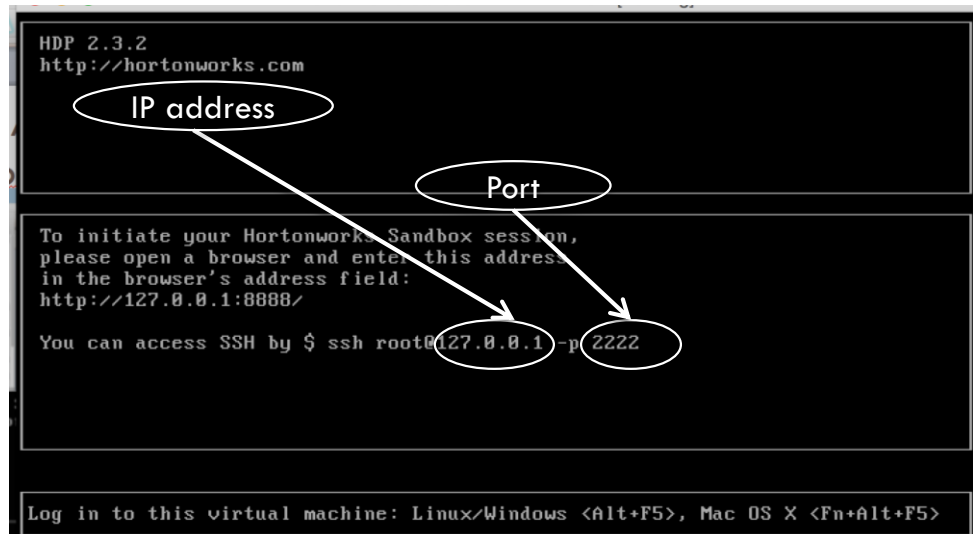


Click here to open File Browser



Choose the HDP_2.3.2_virtualbox.ova file





You could change the RAM value to 4096 MB (4GB instead of 8 GB)

# After Installation
# Hortonworks HDP Sandbox



Double-click

It will lead to the following screen

```
HDP 2.3.2
http://hortonworks.com

        IP address

                          Port

To initiate your Hortonworks Sandbox session,
please open a browser and enter this address
in the browser's address field:
http://127.0.0.1:8888/

You can access SSH by $ ssh root@127.0.0.1 -p 2222


Log in to this virtual machine: Linux/Windows <Alt+F5>, Mac OS X <Fn+Alt+F5>
```

# SSH into the loaded Virtual Machine WINDOWS

- http://www.chiark.greenend.org.uk/~sgtatham/putty/download.html

# SSH into the loaded Virtual Machine MAC

- Open Terminal
- Type "ssh root@127.0.0.1 –p 2222

# Log On and Test HDFS

1. Login Info
   - Username: root
   - Password: hadoop
2. You will be asked to change your password after you login
3. Test HDFS

```
Hortonworks Sandbox with HDP 2.4 [Running]
[root@sandbox ~]# hadoop dfs -ls /
DEPRECATED: Use of this script to execute hdfs command is deprecated.
Instead use the hdfs command for it.

Found 11 items
drwxrwxrwx   - yarn   hadoop          0 2016-03-14 14:19 /app-logs
drwxr-xr-x   - hdfs   hdfs            0 2016-03-14 14:25 /apps
drwxr-xr-x   - yarn   hadoop          0 2016-03-14 14:19 /ats
drwxr-xr-x   - hdfs   hdfs            0 2016-03-14 14:50 /demo
drwxr-xr-x   - hdfs   hdfs            0 2016-03-14 14:19 /hdp
drwxr-xr-x   - mapred hdfs            0 2016-03-14 14:19 /mapred
drwxrwxrwx   - mapred hadoop          0 2016-03-14 14:19 /mr-history
drwxr-xr-x   - hdfs   hdfs            0 2016-03-14 14:42 /ranger
drwxrwxrwx   - spark  hadoop          0 2016-09-04 23:16 /spark-history
drwxrwxrwx   - hdfs   hdfs            0 2016-03-14 14:31 /tmp
drwxr-xr-x   - hdfs   hdfs            0 2016-03-14 14:33 /user
[root@sandbox ~]# _
```

# Test Hive

1. Test Hive

```
[root@sandbox ~]# sudo -u hdfs hive
SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/usr/hdp/2.3.2.0-2950/hadoop/lib/slf4j-log4j12
-1.7.10.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/usr/hdp/2.3.2.0-2950/spark/lib/spark-assembly
-1.4.1.2.3.2.0-2950-hadoop2.7.1.2.3.2.0-2950.jar!/org/slf4j/impl/StaticLoggerBin
der.class]
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF4J: Actual binding is of type [org.slf4j.impl.Log4jLoggerFactory]
WARNING: Use "yarn jar" to launch YARN applications.
SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/usr/hdp/2.3.2.0-2950/hadoop/lib/slf4j-log4j12
-1.7.10.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/usr/hdp/2.3.2.0-2950/spark/lib/spark-assembly
-1.4.1.2.3.2.0-2950-hadoop2.7.1.2.3.2.0-2950.jar!/org/slf4j/impl/StaticLoggerBin
der.class]
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF4J: Actual binding is of type [org.slf4j.impl.Log4jLoggerFactory]

Logging initialized using configuration in file:/etc/hive/2.3.2.0-2950/0/hive-lo
g4j.properties
```
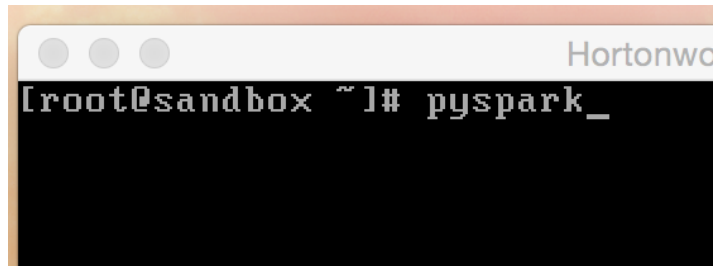
2. You will end up with a screen below

```
Logging initialized using configuration in file:/etc/hive/2.3.2.0-2950/0/hive-lo
g4j.properties
^C[root@sandbox ~]# sudo -u hdfs hive
SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/usr/hdp/2.3.2.0-2950/hadoop/lib/slf4j-log4j12
-1.7.10.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/usr/hdp/2.3.2.0-2950/spark/lib/spark-assembly
-1.4.1.2.3.2.0-2950-hadoop2.7.1.2.3.2.0-2950.jar!/org/slf4j/impl/StaticLoggerBin
der.class]
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF4J: Actual binding is of type [org.slf4j.impl.Log4jLoggerFactory]
WARNING: Use "yarn jar" to launch YARN applications.
SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/usr/hdp/2.3.2.0-2950/hadoop/lib/slf4j-log4j12
-1.7.10.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/usr/hdp/2.3.2.0-2950/spark/lib/spark-assembly
-1.4.1.2.3.2.0-2950-hadoop2.7.1.2.3.2.0-2950.jar!/org/slf4j/impl/StaticLoggerBin
der.class]
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF4J: Actual binding is of type [org.slf4j.impl.Log4jLoggerFactory]

Logging initialized using configuration in file:/etc/hive/2.3.2.0-2950/0/hive-lo
g4j.properties
hive>
```
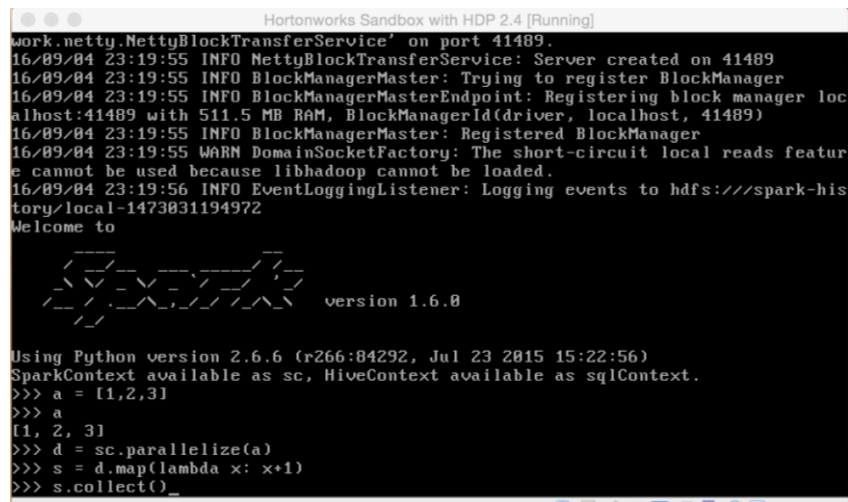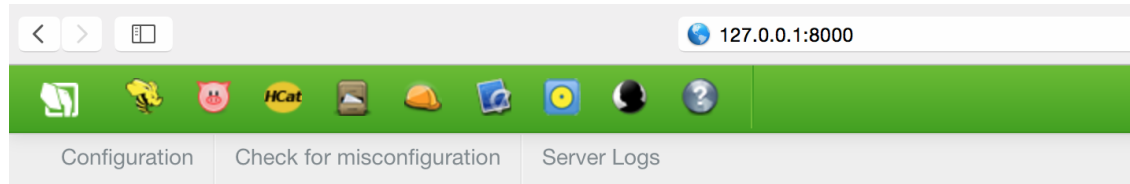
# Test Spark

1. Test Spark
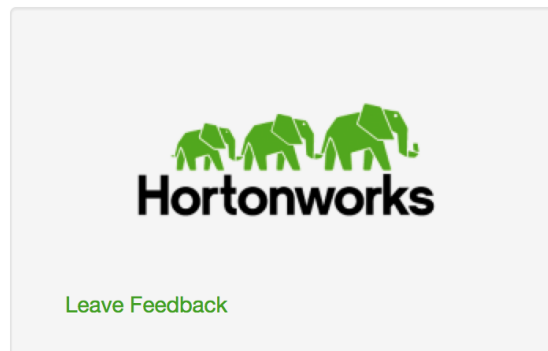


2. You will end up with a screen below

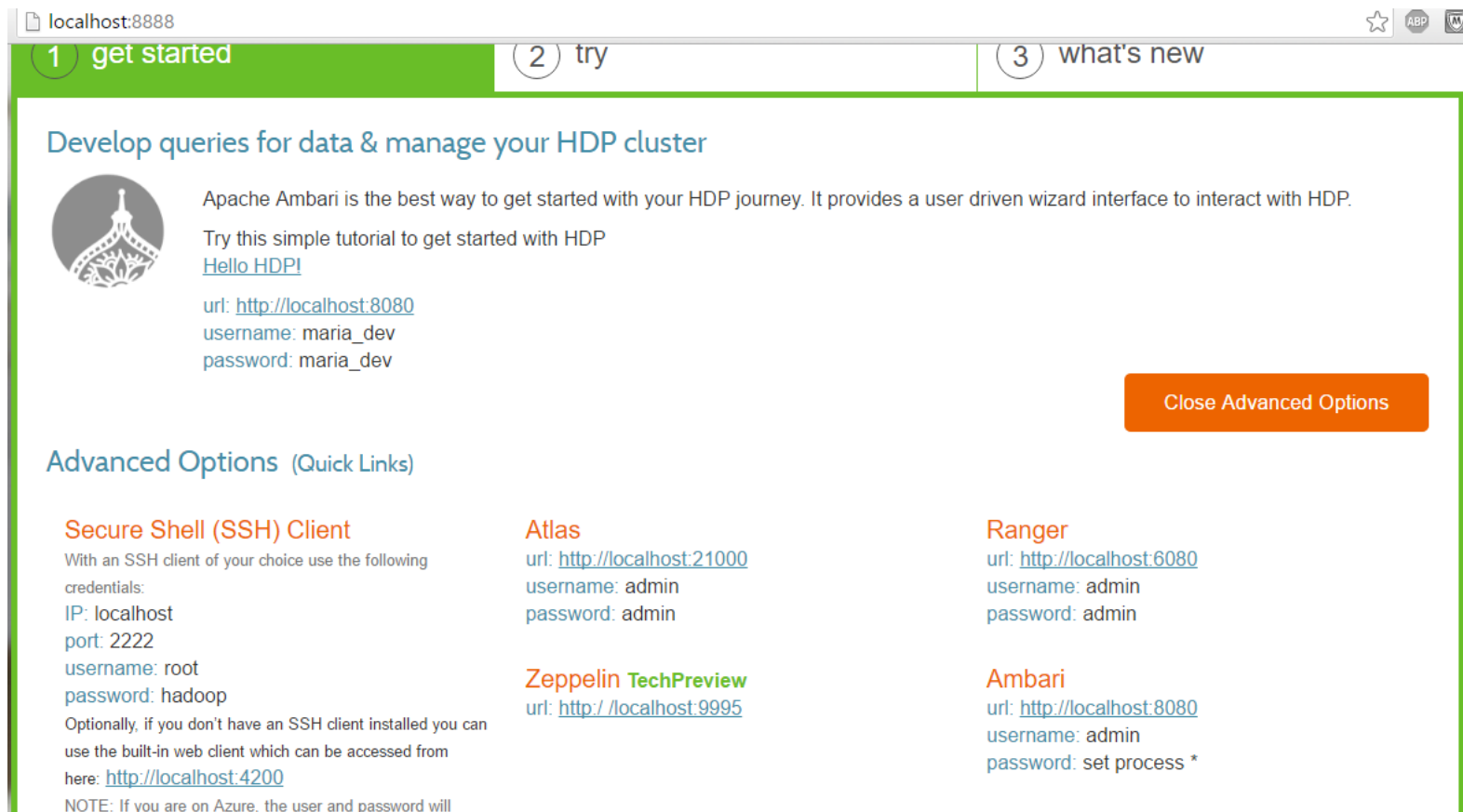# Access Hadoop via Browser Using Hue

## Access Hue via Browser

- http://127.0.0.1:8000

# Welcome Message and Username Information

Go to http://localhost:8888 for welcome screen and Ambari username and password
Click on "Advanced Options" to get SSH and other information

# Browse Through System Setup using Ambari

Go to http://localhost:8080 for Ambari

# Connect to Virtual box using Filezilla

☐ Open Filezilla

☐ Enter Host: sftp://127.0.0.1 or sftp://localhost

☐ Port: 2222

☐ Username & Password – same as your virtual box

☐ Click QuickConnect

# Upload Files to the VirtualBox

- After connecting to the Sandbox access node in Filezilla…
  - The left side shows directories of your local computer ( WINDOWS COMPUTER IF USING LAB MACHINE)
  - The right side box shows directories of your remote machine on Linux
    - In this case the HDP Sandbox (Virtual Box)
- Create a text file called test.txt on your machine with the numbers 1,2,3 written inside
- Upload test.txt to Sandbox
  - On the right-side box, navigate to /root
  - On the left-side box, navigate and find the test.txt file you downloaded
  - Drag and drop "test.txt" into /root/ on the right to the Sandbox

# Linux Command Line

# Try out the following tutorial

- Playing around with big data tools becomes easier with linux command line

- From [http://www.ee.surrey.ac.uk/Teaching/Unix/](http://www.ee.surrey.ac.uk/Teaching/Unix/)

- Try out Tutorial Sections 1, 2, 3, 4, 5 (Section 5.5), and 6

# Python

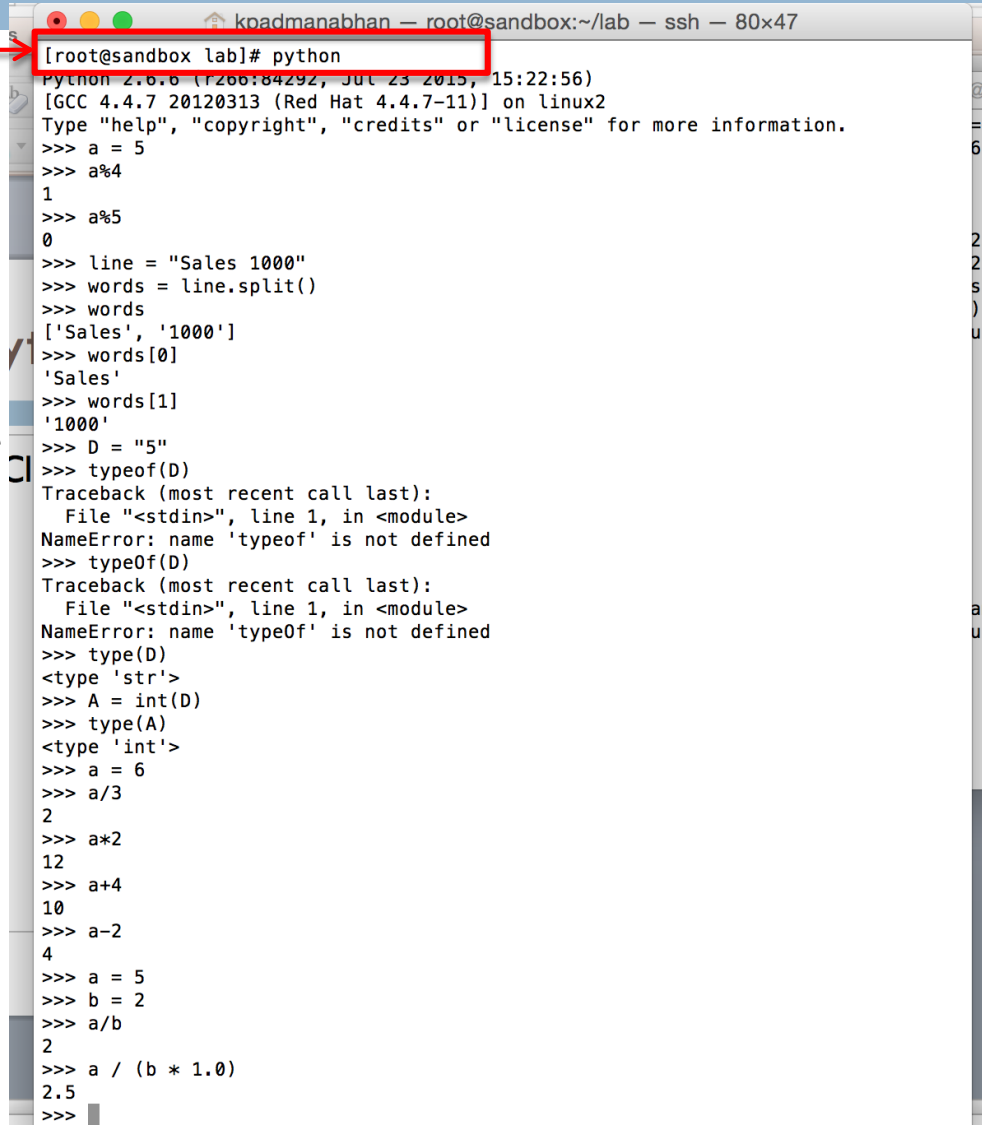# Practice Python Commands

"python" command opens
python shell where we can
try out some commands.
Similar to the "R" shell

1.  a = 5; assigning the number 5 to variable a
2.  a%5; is math modulo operator; it will give the
    value of the reminder when a is divided by 5
3.  line = "Sales 1000"; assigns the string to variable
4.  line.split(); Splits the string into multiple strings;
    Uses "space" to decide where to split
5.  words = line.split(); splits the string and assigns to
    words
6.  A = int(D); Convert string "5" to number 5
7.  type(D); outputs type of D; "string" or "int"
    (integer)
8.  +, -, *, / - same mathematical operators
9.  Notice difference between a/b and a/ (b * 1.0)

http://thepythonguru.com/getting-started-with-
python/
http://www.afterhoursprogramming.com/tutorial
/Python/Introduction/

```
[root@sandbox lab]# python
Python 2.6.6 (r266:84292, Jul 23 2015, 15:22:56)
[GCC 4.4.7 20120313 (Red Hat 4.4.7-11)] on linux2
Type "help", "copyright", "credits" or "license" for more information.
>>> a = 5
>>> a%4
1
>>> a%5
0
>>> line = "Sales 1000"
>>> words = line.split()
>>> words
['Sales', '1000']
>>> words[0]
'Sales'
>>> words[1]
'1000'
>>> D = "5"
>>> typeof(D)
Traceback (most recent call last):
  File "<stdin>", line 1, in <module>
NameError: name 'typeof' is not defined
>>> typeOf(D)
Traceback (most recent call last):
  File "<stdin>", line 1, in <module>
NameError: name 'typeOf' is not defined
>>> type(D)
<type 'str'>
>>> A = int(D)
>>> type(A)
<type 'int'>
>>> a = 6
>>> a/3
2
>>> a*2
12
>>> a+4
10
>>> a-2
4
>>> a = 5
>>> b = 2
>>> a/b
2
>>> a / (b * 1.0)
2.5
>>>
```

# Writing a python script

☐ Copy the following piece of code into *test.py (Keep track of indentation)*

```
import sys
import math
def returnSquare (x):
    return x**2

def main(a):
    print returnSquare(int(a))

if __name__ == "__main__":
    if len(sys.argv) >= 2:
        try:
            main(sys.argv[1])
        except:
            print "Not an Integer"
```

☐ *Execute: python test.py 2*

☐ *Output: 4*

# Try yourself

- Write a python script that will read file shakespere_100.txt and print first 100 lines

- Write a python script that can take a as input from command line and print "even" if even and "odd" if odd

- Write a python script that can take two numbers a & b as input, if a is less than b then calculate a divided by b and if a > b calculate a * b , and print the result rounded to 2 digits