# R·I·T Information Sciences and Technology Department

## NSSA-220 Task Automation Using Interpretive Languages
## Lab 3: Python Data Processing

### INSTRUCTIONS

Complete the tasks in Activity 1 and provide screenshots of the terminal that include your relevant output to prove that you were able to perform each task. In addition, you need to submit your script that generated the output. The lab should be completed and submitted on an individual basis, but feel free to work with other classmates and ask for help from your instructor and TA as needed. When complete, submit the lab document and scripts in a zip file called Lab3.zip to the Lab 3 dropbox. The exact due date will be posted on myCourses.

### PREPARATION

- Read through this document
- Have your Python notes handy

### ACTIVITY SUMMARY

**Activity 0** – Modify your shell prompt
**Activity 1** – Python data processing

### ACTIVITIES

### Activity 0 – Modify your shell prompt

Before you begin working on the lab, you need to modify the prompt in your shell (terminal) to reflect your username, rather than the generic prompt of "student@localhost". To modify your shell prompt, open a Terminal and type "nano .bashrc". We're now going to edit a file that executes every time you start up a terminal in your current session.

```
  GNU nano 2.3.1              File: .bashrc                    Modified

# .bashrc

# Source global definitions
if [ -f /etc/bashrc ]; then
        . /etc/bashrc
fi

# Uncomment the following line if you don't like systemctl's auto-paging featur$
# export SYSTEMD_PAGER=

# User specific aliases and functions
PS1='[username \w]\$ '
```

Now add a line to the end of the file similar to the one above, except that you should use your RIT username in place of the word "username". Save the file and exit nano. Then type exit to leave the Terminal. Open another Terminal and your prompt should now look like the following:

```
[username ~]$
```

We will be looking for the correct username in your screen shots when the lab is graded, so make sure you follow these steps before working on any lab.

## Activity 1 – Python data processing

Perform the following tasks. For each task, include screenshot(s) that clearly indicate the output that proves that the task was accomplished correctly and submit your script that generated the output. Make sure that *your* username is in the screenshot prompt. If you're unable to perform the task as specified, you may receive partial credit by providing the output you were able to get, as well as the script that generated the output, and explain where you had difficulty. Any text you write should be written in red font.

**Task 1 (50 points): Write a Python script called Lab3_Task1.py that takes a single command line argument to represent the input file name that the script will process. Your script will be processing a file called iris_full.txt, which is a famous classification algorithm benchmark data set used in the field of data mining/machine learning.**

**The Iris data set contains 150 records of iris flowers with five fields:**

1. **Sepal length (in cm) of the iris**
2. **Sepal width (in cm) of the iris**
3. **Petal length (in cm) of the iris**
4. **Petal width (in cm) of the iris**
5. **Type of iris (either Iris Setosa, Iris Virginica, or Iris Versicolor)**

**Your script will compute and output a report to the console that includes the following values:**

1. **Minimum value of each of the first four fields**
2. **Maximum value of each of the first four fields**
3. **Average value of each of the first four fields**
4. **The number of occurrences of each type of iris**

**The output should be in the following format:**

Sepal Length: min = <min_value>, max = <max_value>, average = <average_value>

Sepal Width: min = <min_value>, max = <max_value>, average = <average_value>

Petal Length: min = <min_value>, max = <max_value>, average = <average_value>

Petal Width: min = <min_value>, max = <max_value>, average = <average_value>

Iris Types: Iris Setosa = <num_iris_setosa>, Iris Versicolor = <num_iris_versicolor>, Iris Virginica = <num_iris_virginica>

**The correct answers for all of these outputs are shown in the comments section of iris_full.txt**

<insert screenshot here>

**Your script should contain the following functions:**

1. **read_data: this function takes two arguments – a file name to read and a List L to store all the records from the file in. Your script must read the file *as is*. Do not modify the file prior to reading it.**

2. **process_numeric_field**: the function takes two arguments – the List of records L and the field number to process (1, 2, 3, or 4). The function returns three values – min, max, and average of the field.
3. **count_iris_types**: the function takes one argument – the List of records L. The function returns three values – the number of Iris Setosa, Iris Versicolor, and Iris Virginica in the data set.

Extra credit (10 points): process_numeric_field also computes the standard deviation of a field. The correct answers for each field are also shown in iris_full.txt (refer to SD). If you compute the standard deviation of a field, output it along with the rest of the statistics for a field by including ", standard deviation = <std_value>" in the output.

For this task, submit the Lab3_Task1.py script as part of your Lab3.zip file.

Task 2 (50 points): Write a Python script called Lab3_Task2.py that takes two command line arguments, one for each of the files that the script will process.

Your company's security department recently discovered a network breach and they suspect that the CentOS VM image you have been using was compromised. Specifically, they're worried that the executables in /usr/bin could have been modified with versions of those same executables that now contain malicious code.

You have been provided with a file called md5_new.txt that contains the MD5 message digests for nearly all of the executables in /usr/bin. MD5 message digests are 128-bit hashes that have been computed for the executables shown in md5_new.txt. Each line in the file contains the name of an executable in /usr/bin and its MD5 message digest that was computed before the network breach. These message digests are important because they should be unique to those specific executables and if the message digest is now different, that executable was affected by the breach. MD5 message digests for any file may be computed using a system call to md5sum.

In addition, you've been provided with a file called md5_original.txt, which contains the MD5 message digests for the executables in /usr/bin prior to the suspected network breach.

Your script will output the names of executables that may have been compromised to the console using the following format:

<affected_executable_name_1>: MD5 original = <md5_from_file>, MD5 new = <md5_from_your_script>

<affected_executable_name_2>: MD5 original = <md5_from_file>, MD5 new = <md5_from_your_script>

…

<affected_executable_name_n>: MD5 original = <md5_from_file>, MD5 new = <md5_from_your_script>

<insert screenshot here>

The only requirements for your code are that the script accepts the file names as command line arguments and it compares the message digests between md5_new.txt and md5_original.txt.

For this task, submit the Lab3_Task2.py script as part of your Lab3.zip file.