

# Using R on the HPC

Alex Townsend



# Introduction to R

- R is a statistical analysis tool and a programming language.
- R is widely used for data analysis in a variety of fields.
- R is free and open source!

# Using R on the HPC

- R is available on both the Spear and HPC Systems.
  - There are two versions of R available.
    - The default version can be loaded by typing the following into the terminal
      - `R`
    - The other version, which is the most up-to-date and is recommended for use can be loaded with the following commands:
      - `module load R`
      - `R`

# Using R on the HPC

- You can run R on the Login Nodes
  - This limits you to one node and up to 24 cores.
- You can also submit your R job to the Compute Nodes.
  - This allows you to use multiple nodes, each of which have a large number of cores, usually between 24 and 48.
  - This also allows you to use GPUs.

# The Advantages of R on HPC

- R can be run normally on any laptop. Why use HPC?
  - Parallel Computing!
    - The parallel package ~ Multicore
    - The rmpi package ~ Distributed
  - GPU Computing!
    - gpuR and tensorflow Packages

# Today's Example

- For the example today, we will be tuning a simple R script using a built-in dataset of SAT and ACT scores. This will go in 2 steps.
  - Step 1: Write a simple script for Serial Computation
  - Step 2: Modify the script for Multicore Computation

# Step 1

```
1 # load the Package with the Data
2 library(psych)
3
4 # Get the SAT and ACT Score Dataset
5 dataset <- sat.act
6
7 # Clean the Dataset (Remove Missing Data)
8 cleandata <- na.omit(dataset)
9
0 # Now get the Column Means for the data
1 means <- colMeans(cleandata, na.rm=TRUE)
2
3 # Now compute the Column-Wise Variances for the data
4 vars <- apply(cleandata, 2, var)
5
```

# Step 1

```
1 #!/bin/bash
2
3 #SBATCH --job-name="R Template"
4
5
6
7 #SBATCH --mail-type=ALL
8
9 #SBATCH -n 24
10
11
12 #SBATCH -p genacc_q
13
14 #SBATCH -t 14-00:00:00
15
16
17
18
19
20
21
22 module load R
23
24
25
26 R CMD BATCH parallel_script.R
```



# Step 2

```
1 # load the Package with the Data
2 library(psych)
3
4 # Load the PARALLEL library for Multicore Computing!
5 library(parallel)
6
7 # Get the SAT and ACT Score Dataset
8 dataset <- sat.act
9
10 # Clean the Dataset (Remove Missing Data)
11 cleandata <- na.omit(dataset)
12
13 # Determine how many cores you have available
14 cores <- detectCores() - 1
15 print(cores)
16
17 # Build a virtual "cluster" out of these cores
18 cluster <- makeCluster(cores)
19
20 # Now get the Column Means for the data
21 means <- parLapply(cluster, cleandata, mean)
22 print(means)
23
24 # Now compute the Column-Wise Variances for the data
25 vars <- parLapply(cluster, cleandata, var)
26 print(vars)
27
28 # Now stop the cluster and release the resources
29 stopCluster(cluster)
```

# Step 2

```
1 #!/bin/bash
2
3 #SBATCH --job-name="R Template"
4
5
6
7 #SBATCH --mail-type=ALL
8
9 #SBATCH -n 24
0
1
2 #SBATCH -p genacc_q
3
4 #SBATCH -t 14-00:00:00
5
6
7
8
9
0
1
2 module load R
3
4
5
6 R CMD BATCH parallel_script.R
```

# Advanced Topics

- R also has packages available to make that same script useful for Distributed Computation
  - For this, we would use the Rmpi package.

# Advanced Topics

- R also has packages available to make that same script useful for GPU Computation
  - gpuR and tensorflow