

# NISS/CANSSI WORKSHOP

NISS

# NISS

NISS/CANSSI WORKSHOP

R & SPARK: TOOLS FOR DATA  
SCIENCE WORKFLOWS  
APRIL 12-13, 2018



## R & SPARK: TOOLS FOR DATA SCIENCE WORKFLOWS

**CLASS CAPACITY:** 25

**FEES:**

- 760 USD or 985 CAD for students, or employees of NISS Affiliates
- 990 USD or 1,285 CAD for non-NISS Affiliates

**COURSE OUTLINE:** R is a flexible, extensible statistical computing environment, but it is limited to single-core execution. Spark is a distributed computing environment which treats R as a first-class programming language. This course introduces data structures in R and their use in functional programming workflows relevant to data science.

The course covers the initial steps in the data science process:

- Extracting data from source systems
- Transforming data into tidy form
- loading data into distributed file systems, distributed data warehouses, and NoSQL databases, i.e., ETL.

This workflow is illustrated by using the SparkR and sparklyr package frontends to Spark from R.

SparkR and sparklyr are then used as interfaces for modeling big data using regression and classification supervised learning methods. Unsupervised learning methods, such as clustering and dimension reduction, are also covered. Additional methods, such as gradient boosting and deep learning, are illustrated using the h2o and rsparkling R packages. Finally, methods for analyzing streaming data are presented. The course finishes with an in-depth example. The infrastructure and content is containerized for easy download to your laptop using Docker.

**PREREQUISITES FOR THIS COURSE:** Differential calculus, basic matrix algebra, a statistics course covering regression, basic R.

**OPERATING SYSTEMS:** MacOS 10.11 (El Capitan) or higher or Windows 10 Professional. Students must bring their own laptops.

**REGISTER:** Register online with a credit card at <https://www.niss.org/events/r-and-spark-tools-data-science-workflows-2>

You can also call (202) 862-4316 or write to [officeadmin@NISS.org](mailto:officeadmin@NISS.org) to register

**CONTACT US:** Direct questions about this course to the Instructor E. James Harner at [eharner@mail.wvu.edu](mailto:eharner@mail.wvu.edu) or call him on his cell phone at **304-376-4170**.

For other questions, contact [officeadmin@NISS.org](mailto:officeadmin@NISS.org)



**INSTRUCTOR:**  
**E. JAMES  
HARNER**

E. James Harner is Professor Emeritus of Statistics at West Virginia University (WVU). He was the Chair of the Department of Statistics for 17 years and the Director of the Cancer Center Bioinformatics Core for 15 years at WVU. Currently, he is the Chairman of the Interface Foundation of North America which has partnered with the American Statistical Association to organize the annual Symposium on Data Science and Statistics (SDSS) beginning in May, 2018. The areas of his technical and research expertise include: bioinformatics, high-dimensional modeling, high-performance computing, streaming and big data modeling and statistical machine learning.

**National Institute of  
Statistical Sciences**

1150 Connecticut  
Avenue NW, 9th Floor,  
Washington, DC 20036;  
**Tel:** (202) 862-4316;  
**Fax:** (202) 828-4130