# IIT School of Applied Technology
## ILLINOIS INSTITUTE OF TECHNOLOGY
**information technology & management**

# 527 Data Analytics

March 22,24 2016

Week 11 Presentation

# Week 11 Topic: Agenda

**ITM - 527**

- Week 11 Assignment Review
- Hadoop in more detail
- Hadoop and R in SAS

Big Data:
- What is Big Data?
- What is Hadoop?
- Big Data Analytics
- Intro to R
- Installing R or RStudio

# Week 11 Topic:
# Assignment 1 – Cluster Analysis in R

(Worth 10 points) Due March 31st

◆ Read Chapter 10 of Data Smart

◆ Follow directions as stated and perform cluster analysis for k=3 in R

◆ Submit code and results in a word document.

ITM – 527

# Week 11 Topic:
# Assignment 2 – Exec Comp Review

ITM - 527

(Worth 10 points) Due March 31st

◆ Read: https://pdfs.semanticscholar.org/52c6/d51d85e9a2c529d696cb83b6ad1d7daa45ca.pdf

◆ Download/Collect DEF 14A filings for 5 companies to analyze from www.sec.gov

◆ Summarize, for each company, its CEO name, Salary, Vested Stock Awards, and Option Awards for all available years in the filing.

◆ Determine total vested stock ownership of the CEO and salary trend for all available years in the filing.

◆ Describe (in numbered steps) how to automate the above task in SAS or R.

ITM - 527

# Week 11 Topic:
# SAS and Hadoop

◆ The Apache Software Foundation[1] provides support for open-source software projects. One such project is Apache Hadoop[2].

◆ *Hadoop* is a framework for distributed processing of large data sets across a cluster of computers using simple programming models[2]

[1] www.apache.org
[2] hadoop.apache.org

A cluster of computers:



**5**

# Week 11 Topic:
# Cluster of Computers?

http://en.wikipedia.org/wiki/Computer_cluster

## Computer cluster

From Wikipedia, the free encyclopedia

*Not to be confused with data cluster or computer lab.*

A **computer cluster** consists of a set of loosely connected or tightly connected computers that work together so that in many respects they can be viewed as a single system.

The components of a cluster are usually connected to each other through fast local area networks ("LAN"), with each *node* (computer used as a server) running its own instance of an operating system. Computer clusters emerged as a result of convergence of a number of computing trends including the availability of low cost microprocessors, high speed networks, and software for high performance distributed computing.

Technicians working on a large Linux cluster at the Chemnitz University of Technology, Germany

ITM - 527

**6**

# Week 11 Topic:
# Revisit of Core Hadoop Modules

◆ Core Hadoop modules include:

| | |
|---|---|
| **HDFS (Hadoop Distributed File System)** | a file system that distributes large files across the Hadoop cluster of computers |
| **Hadoop YARN** | a framework for job scheduling and cluster resource management |
| **Hadoop MapReduce** | a YARN-based system for parallel processing of large data sets |

◆ These modules automate the process of reading, writing, and processing large files in a distributed environment, freeing programmers to write programs to process the data as if they were using a single computer.

**ITM - 527**

**7**

# Week 11 Topic:
# File management in Hadoop

◆ HDFS is hierarchical with LINUX style paths and file ownership and permissions.
- HADOOP FS commands are similar to LINUX commands.
- HDFS in not built into the operating system.
- Files are append-only after they are written.

HADOOP FS commands from LINUX command prompt:

```
$ hadoop fs -ls /user/student
Found 4 items
drwxr-xr-x  - student1 sasapp 0 2014-05-30 20:00 /user/student1/.Trash
drwx------  - student1 sasapp 0 2014-05-30 10:05 /user/student1/.stage
drwxr-xr-x  - student1 sasapp 0 2014-05-28 15:25 /user/student1/data
drwxr-xr-x  - student1 sasapp 0 2014-05-28 13:59 /user/student1/users
$ hadoop fs -mkdir /user/student1/newdir
$
```
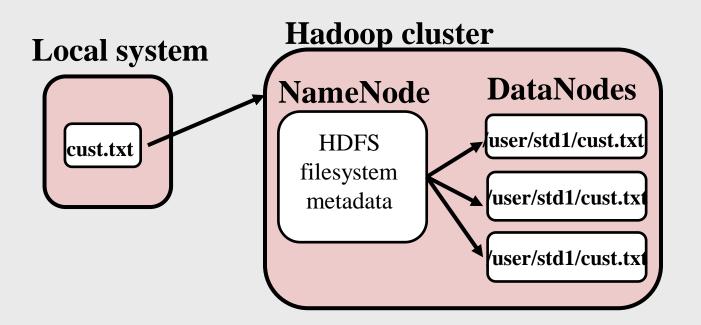
ITM - 527

# Week 11 Topic:
# File management in Hadoop

◆ The following HDFS command moves a local file into the HDFS cluster:

```
$ hadoop fs copyfromlocal="cust.txt"  out="/user/std1"
```

◆ An HDFS *NameNode* on a root node machine distributes the file to each of the *DataNode* machines, and provides access to the distributed file.

**Local system**

**Hadoop cluster**

cust.txt

**NameNode**

HDFS filesystem metadata

**DataNodes**

/user/std1/cust.txt

/user/std1/cust.txt

/user/std1/cust.txt

**ITM - 527**

**9**

# Week 11 Topic:
# Characteristics of Hadoop

**ITM - 527**

- Open-source
- Simple to use distributed file system
- Supports highly parallel processing
- It's scalable, so it's suitable for massive amounts of data
- It is designed to work on low-cost hardware
- It's fault tolerant at the data level
  - automatic replication of data
  - automatic fail-over

# Week 11 Topic:
# Hadoop Distribution

ITM - 527

◆ Apache Hadoop is open-source technology developed by the Apache Software Foundation

- a community of developers and users
- a non-profit

◆ Several commercial vendors provide, support, and augment Apache Hadoop with their own distribution. These include:

- Cloudera
- Hortonworks
- BigInsights

11

# Week 11 Topic: SAS Hadoop Interface

ITM - 527

| Tool | Purpose | Product |
|------|---------|---------|
| FILENAME statement | Allows the DATA step to read and write HDFS data files. | Base SAS |
| PROC HADOOP | Copy or move files between SAS and Hadoop. Execute MapReduce and Pig code in Hadoop. Execute Hadoop file system commands to manage files and directories. | Base SAS |
| SQL Pass-Through | Submit HiveQL queries and other HiveQL statements from SAS directly to Hive for Hive processing. Query results are returned to SAS. | SAS/ACCESS Interface to Hadoop |
| LIBNAME Statement For Hadoop | Access Hive tables such as SAS data sets using the SAS programming language. SAS/ACCESS engine translate SAS language into HiveQL and attempts do convert as much of the processing into HiveQL as possible before returning results to SAS. | SAS/ACCESS Interface to Hadoop |

# Week 11 Topic:
# Base SAS Tools for Hadoop

ITM - 527

◆    The Hadoop FILENAME statement enables you to do the following:
- upload local data to Hadoop using the DATA step
- read data from Hadoop using the DATA step

◆    PROC HADOOP enables you to do the following:
- submit Hadoop file system (HDFS) commands
- submit MapReduce programs
- submit Pig language code

# Week 11 Topic:
# Hadoop Config.xml File

◆    The FILENAME statement for Hadoop and PROC HADOOP require an option that specifies a Hadoop configuration file.

- The configuration file defines how to connect to Hadoop and other Hadoop system information.
- The file must be accessible to the SAS client application.
- A SAS administrator commonly manages this configuration for the SAS users.
- This file is often referred to as the Hadoop core-site.xml file.

ITM – 527

# Week 11 Topic:
# Hadoop Config.xml File

```xml
<?xml version="1.0" encoding="UTF-8"?>
- <configuration>
  - <property>
        <name>fs.default.name</name>
        <value>hdfs://sasserver.demo.sas.com:8020</value>
    </property>
  - <property>
        <name>mapred.job.tracker</name>
        <value>sasserver.demo.sas.com:8021</value>
    </property>
  - <property>
        <name>dfs.http.address</name>
        <value>sasserver.demo.sas.com:50070</value>
    </property>
  - <property>
        <name>dfs.secondary.http.address</name>
        <value>sasserver.demo.sas.com:50090</value>
    </property>
  - <property>
        <name>mapred.job.tracker.http.address</name>
        <value>sasserver.demo.sas.com:50030</value>
    </property>
</configuration>
```
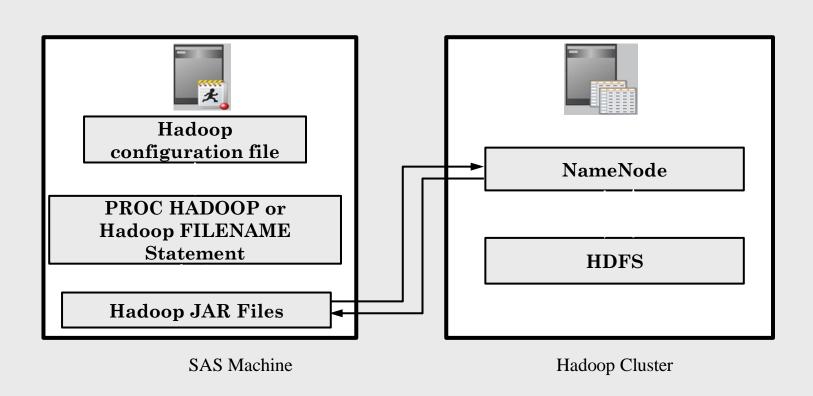
ITM - 527

**15**

# Week 11 Topic:
# Hadoop JAR File

**ITM – 527**

◆ A collection of Hadoop JAR files is also required on the SAS client machine.

- An environment variable SAS_HADOOP_JAR_PATH on the SAS client machine defines location of the Hadoop JAR files

- The Hadoop JAR files must be compatible with the specific Hadoop implementation and can be copied from the Hadoop server.

- A Hadoop system administrator commonly manages the configuration of the Hadoop JAR files for the SAS users.

# Week 11 Topic:
# Hadoop Interface to SAS

**ITM - 527**



| Hadoop configuration file |
|---|

| PROC HADOOP or Hadoop FILENAME Statement |
|---|

| Hadoop JAR Files |
|---|

| NameNode |
|---|

| HDFS |
|---|

SAS Machine

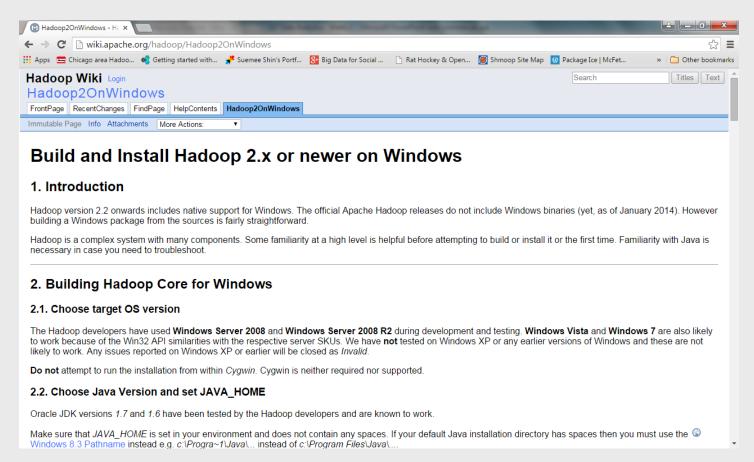Hadoop Cluster

# Week 11 Topic:
# Installing Hadoop

ITM - 527

https://hadoop.apache.org/docs/current/hadoop-project-dist/hadoop-common/SingleCluster.html (using a Linux emulator)

# Week 11 Topic:
# Installing Hadoop

http://wiki.apache.org/hadoop/Hadoop2OnWindows

ITM - 527

## Build and Install Hadoop 2.x or newer on Windows

### 1. Introduction

Hadoop version 2.2 onwards includes native support for Windows. The official Apache Hadoop releases do not include Windows binaries (yet, as of January 2014). However building a Windows package from the sources is fairly straightforward.

Hadoop is a complex system with many components. Some familiarity at a high level is helpful before attempting to build or install it or the first time. Familiarity with Java is necessary in case you need to troubleshoot.

### 2. Building Hadoop Core for Windows

#### 2.1. Choose target OS version

The Hadoop developers have used **Windows Server 2008** and **Windows Server 2008 R2** during development and testing. **Windows Vista** and **Windows 7** are also likely to work because of the Win32 API similarities with the respective server SKUs. We have **not** tested on Windows XP or any earlier versions of Windows and these are not likely to work. Any issues reported on Windows XP or earlier will be closed as *Invalid*.

**Do not** attempt to run the installation from within *Cygwin*. Cygwin is neither required nor supported.

#### 2.2. Choose Java Version and set JAVA_HOME

Oracle JDK versions *1.7* and *1.6* have been tested by the Hadoop developers and are known to work.

Make sure that *JAVA_HOME* is set in your environment and does not contain any spaces. If your default Java installation directory has spaces then you must use the
Windows 8.3 Pathname instead e.g. *c:\Progra~1\Java\...* instead of *c:\Program Files\Java\...*

# Week 11 Topic: MapReduce

**ITM - 527**

◆ *MapReduce* is a framework written in Java that is built into Hadoop. It automates the distributed processing of data files.

| | |
|---|---|
| *map* | processing of individual rows (filtering, row calculations) |
| *shuffle and sort* | grouping rows for summarization |
| *reduce* | summary calculations within groups |

◆ The MapReduce framework coordinates multiple mapping, sorting, and reducing tasks that execute in parallel across the computer cluster.

**20**

**ITM - 527**

# Week 11 Topic:
# Executing MapReduce

A JAVA programmer can write code that calls the following 'functions' that are built into the MapReduce framework:

- an input reader
- a Map function
- a partition function
- a compare function
- a Reduce function
- an output writer

More recently, Pig and Hive projects have been developed for Apache Hadoop. These provide a higher-level language interface that compiles to MapReduce.

# Week 11 Topic:
# Pig and Hive

Pig and Hive provide less complex higher-level programming methods for parallel processing of Hadoop data files.

| | |
|---|---|
| *Pig* | A platform for data analysis that includes stepwise procedural programming that converts to MapReduce. |
| *Hive* | A data warehousing framework to query and manage large data sets stored in Hadoop. Provides a mechanism to structure the data and query the data using an SQL-like language called HiveQL. Most HiveQL queries are compiled into MapReduce programs. |

ITM – 527

# Week 11 Topic:
# SAS and R

http://www.sas.com/content/dam/SAS/en_us/doc/conclusionpaper1/use-of-open-source-is-growing-107429.pdf:

SAS/IML – a matrix manipulation product that supports matrix-vector computation – gives users the flexibility to create custom functions. It also:

◆ Takes advantage of built-in functions, subroutines and SAS procedures.
◆ Enables interfaces with R language so users can submit R code within SAS.
◆ Moves between SAS and R data structures.
◆ Enables R functions and packages.

YouTube on how to do this:

◆ www.youtube.com/watch?v=rUaTTre24kI
◆ www.youtube.com/watch?v=nmRQ3MtkG6A

**ITM - 527**

# Week 11 Topic:
# R Reference and Exercises

ITM - 527

SAS and R Reference:

http://sas-and-r.blogspot.com/

Sample Chapters and Exercises:

http://www3.amherst.edu/~nhorton/sasr2/

# Week 11 Topic:
# Profiling & Tax Strategy Submission

(Worth 20 points)

◆ Fill in the Profile Matrix Spreadsheet following the example in class.

◆ Determine your own Tax Groups. Similar to the example in class. Provide answers in the same Matrix spreadsheet.

◆ Determine the tax strategy per Tax Group(s) and provide the overall tax increase amount. Provide answers in the same Matrix workbook.

ITM - 527

# Week 11 Topic:
# What is Big Data?

ITM - 527

*Very large, distributed aggregations of loosely structured data – often incomplete and inaccessible*

**Before:**

◆ Used to mainly refer to data on web behavior or social network interactions. Even 5 years ago was something to think about but nothing to set aside significant IT budget on if you were a non-internet organization.

◆ Used to also be a specialized capability that companies like comScore, Yahoo, Google, and Amazon developed to collect/analyze internet browsing behavior and purchase data on consumers.

◆ A significant application of it's technology is on fraud detection and surveillance activities.

**Now:**

◆ It has broadened to include collection and analysis of any unstructured data that are core to businesses including traditional organizations such as hospitals, banks, and supply chains.

◆ These traditional companies used to ETL data from systems and paper based functions into databases for analytical work. This process took a lot of time and effort. Now, the idea with Big Data is that we EL or L first then ET or T as needed for analytical requests on the fly.

**26**

# Week 11 Topic:
# Big Data / Traditional DW

ITM - 527

**Traditional Data Warehouse:**

◆ Records from internal transaction systems or subscription based data provider

◆ Data is centralized and uniform

◆ Batch updates of new data, mostly historical data

◆ Analytics designed against stable (data model) environment

◆ Scheduled production reports

→ Key is that the data is known and data quality is part of the processing. Long lead time in data availability but less processing to query the available data

**Big Data Analytical Environments:**

◆ Data available from various sources inside and outside of organizations, including DWs

◆ Distributed data in various formats

◆ Iteration needed to model/test data

◆ May require large amount of memory/resources, on the fly, to process request/query

→ Key is that all data is available in its native form at any given time. Short lead time in data availability but longer processing time to query the available data if needing additional normalization steps

# Week 11 Topic:
# What is Hadoop?

*Open-source software that enables reliable, scalable, distributed computing on clusters of inexpensive servers*

◆ Reliable: The software is fault tolerant, it expects and handles hardware and software failures

◆ Scalable: Designed for massive scale of processors, memory, and local attached storage

◆ Distributed: Handles replication. Offers massively parallel programming model, MapReduce

**Hadoop Framework:**

**Non-Relational DB**
Fine-grained data handling

**Hive**
Data warehouse that provides SQL interface. Data structure is projected ad hoc onto unstructured underlying data

**HBase**
Column oriented, schema-less, distributed database modeled after Google's BigTable. Random realtime read/write

**Scripting**

**Pig**
Platform for manipulating and analyzing large data sets. Scripting language for analysts

**Machine Learning**

**Mahout**
Machine learning libraries for recommendations, clustering, classfication and itemsets

**MapReduce**
- Parallel programming
- Large block data handling (e.g. 64MB)

**Hadoop Common**

**HDFS**
Distributes & replicates data across machines

**MapReduce**
Distributes & monitors tasks, restarts failed work

28

# Week 11 Topic:
# Hadoop Framework

*Hadoop is designed to process terabytes and even petabytes of unstructured and structured data. It breaks large workloads into smaller data blocks that are distributed across a cluster of commodity hardware for faster processing. And it's part of a larger framework of related technologies:*

◆ HDFS: Hadoop Distributed File System

◆ HBase: Column oriented, non-relational, schema-less, distributed database modeled after Google's BigTable. Promises "Random, real-time read/write access to Big Data"

◆ Hive: Data warehouse system that provides SQL interface. Data structure can be projected ad hoc onto unstructured underlying data

◆ Pig: A platform for manipulating and analyzing large data sets. High level language for analysts

◆ ZooKeeper: a centralized service for maintaining configuration information, naming, providing distributed synchronization, and providing group services. A tool for configuring and synchronizing Hadoop clusters.

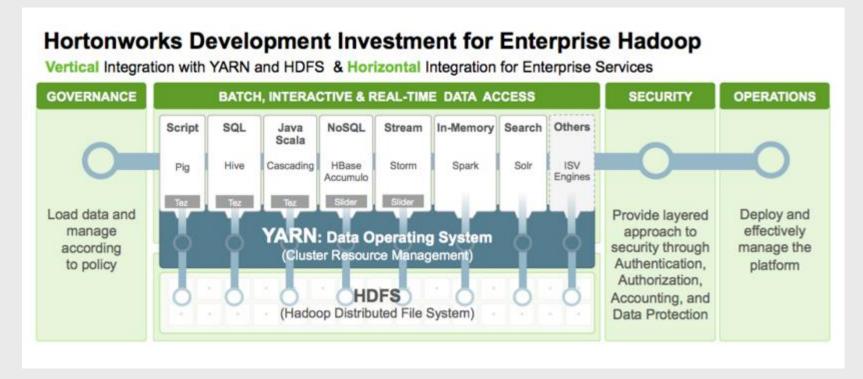ITM - 527

**29**

# Week 11 Topic:
# When to use Hadoop?

◆ Complex but parallelizable algorithms needed, such as geo-spatial analysis or genome sequencing

- Complex information processing is needed
- Unstructured data needs to be turned into structured data
- Queries can't be reasonably expressed using SQL
- Heavily recursive algorithms e.g., need for parallel processing like machine learning

◆ Data sets are too large to fit into database RAM, discs, or require too many cores (10's of TB up to PB). Fault tolerance and scale is a challenge.

◆ Data value does not justify expense of constant real-time availability, such as archives or special interest info, which can be moved to Hadoop and remain available at lower cost and lower latency. Results are not needed in real time

ITM – 527

# Week 11 Topic: Blueprint for Hadoop

ITM - 527

http://hortonworks.com/blog/announcing-hdp-2-2/

(https://hortonworks.com/wp-content/uploads/2013/10/Big-Data-Reference-Architecture_4AA5-6136ENW.pdf ):



**Hortonworks Development Investment for Enterprise Hadoop**

**Vertical** Integration with YARN and HDFS & **Horizontal** Integration for Enterprise Services

| GOVERNANCE | BATCH, INTERACTIVE & REAL-TIME DATA ACCESS | | | | | | | | SECURITY | OPERATIONS |
|---|---|---|---|---|---|---|---|---|---|---|
| | Script | SQL | Java Scala | NoSQL | Stream | In-Memory | Search | Others | | |
| | Pig | Hive | Cascading | HBase Accumulo | Storm | Spark | Solr | ISV Engines | | |
| | Tez | Tez | Tez | Slider | Slider | | | | | |
| Load data and manage according to policy | **YARN: Data Operating System** (Cluster Resource Management) | | | | | | | | Provide layered approach to security through Authentication, Authorization, Accounting, and Data Protection | Deploy and effectively manage the platform |
| | **HDFS** (Hadoop Distributed File System) | | | | | | | | | |

31

# Week 11 Topic:
# Blueprint for Hadoop

**ITM - 527**

(https://hortonworks.com/wp-content/uploads/2013/10/Big-Data-Reference-Architecture_4AA5-6136ENW.pdf ):

Key highlights of HDP 2.2 include the following:

◆ Batch and interactive SQL queries via Apache Hive and Apache Tez, along with a cost-based optimizer powered by Apache Calcite

◆ High-performance ETL via Pig and Tez

◆ Stream processing via Apache Storm and Apache Kafka

◆ YARN labels

◆ Search via Apache Solr

◆ Streamlined cluster operations via Apache Ambari

◆ Data lifecycle management via Apache Falcon

◆ Perimeter security via Apache Knox

◆ Centralized security administration for HDFS, Hive, HBase, Storm, and Knox via Apache Ranger

# Week 11 Topic:
# DB vs Hadoop

ITM - 527

|  | Database | Hadoop (HDFS+MapReduce) |
|---|---|---|
| Function | - Defined as a database<br>- Built to support structured data | - Defined as a distributed file system – NoSQL (Not only SQL)<br>- Built to support unstructured as well as structured data |
| Schema | - Stored in rows and columns<br>- Required on write | - Stored as key/value pairs<br>- Required on read |
| Processing | - Unlimited read and write<br>- Pooled processing across nodes<br>- Configuration performed during install | - Write once, process on read<br>- Recommend processing per node<br>- Use YARN to manage resources |
| Query | - Returns one answer across all data<br>- Immediate results, real time | - Returns all answers to process on<br>- Requires time to process raw files or can build a database on top of Hadoop (Hbase) |
| API | SQL, Native, ODBC, JDBC | MapReduce |
| Scale | - Scaling is a challenge<br>- Vertical scaling, boost server for performance | - Built to scale with ease<br>- Horizontal scaling, add more nodes |
| TCO | - Typically, licensed software on dedicated servers<br>- Functions built out of the box | - Open source software with commodity servers<br>- Need to build traditionally available functions e.g., security, logging, ETL, server brokering, etc. |
| When to use | Workloads are constant and predictable | Challenged by increasing data demands |

33

# Week 11 Topic:
# Big Data Playbook

**ITM – 527**

This may counter what we just discussed in the previous slide but does offer a good perspective.

http://www.sas.com/content/dam/SAS/en_us/doc/whitepaper1/non-geeks-big-data-playbook-106947.pdf

# Week 11 Topic:
# Big Data Analytics

ITM - 527

◆ Big Data Analytics is going to be different than what was offered via traditional architectures. Data being evaluated is different. Technology involved is different. The massive amount of data also makes you think of different ways to perform analytics on the data

◆ Some highlights of what to expect:

- For batch oriented data analysis / prototyping with lower processing power required use:

  → R, Python, Matlab

- For low latency / real time data analysis / prototyping with lower processing power required use:

  → Spark (Scala)

- For low latency / real time data analysis / prototyping with lots of processing power required use:

  → C/C++, Fortran

- Other incumbents to the low latency, fast processing arena: Julia, Python, Spark

- SAS, SPSS, and other analytic applications have Hadoop APIs that can be leveraged

# Week 11 Topic:
# R

ITM – 527

http://en.wikipedia.org/wiki/R_(programming_language):

◆ *Ross Ihaka and Robert Gentleman* created the open-source language R in 1995 as an implementation of the S programming language. The purpose was to develop a language that focused on delivering a better and more user-friendly way to do data analysis, statistics and graphical models. At first, R was primarily used in academics and research, but lately the enterprise world is discovering R as well. This makes R one of the fastest growing statistical languages in the corporate world.

◆ One of the main strengths of R is its huge community that provides support through mailing lists, user-contributed documentation and a very active Stack Overflow group. There is also CRAN (http://cran.r-project.org/), a huge repository of curated R packages to which users can easily contribute.  These packages are a collection of R functions and data that make it easy to immediately get access to the latest techniques and functionalities without needing to develop everything from scratch yourself.

ITM - 527

# Week 11 Topic:
# Introducing R Project

https://www.r-project.org/:

# Week 11 Topic:
# Introducing CRAN

ITM - 527

http://cran.r-project.org/

# Week 11 Topic:
# R and CRAN

**What are R and CRAN?**

◆ R is 'GNU S', a freely available language and environment for statistical computing and graphics which provides a wide variety of statistical and graphical techniques: linear and nonlinear modelling, statistical tests, time series analysis, classification, clustering, etc. Please consult the R project homepage for further information.

◆ CRAN is a network of ftp and web servers around the world that store identical, up-to-date, versions of code and documentation for R. Please use the CRAN mirror nearest to you to minimize network load.

**Submitting to CRAN**

◆ To "submit" a package to CRAN, check that your submission meets the CRAN Repository Policy and then use the web form.

◆ If this fails, upload to ftp://CRAN.R-project.org/incoming/ and send an email to CRAN@R-project.org following the policy. Please do not attach submissions to emails, because this will clutter up the mailboxes of half a dozen people.

◆ Note that we generally do not accept submissions of precompiled binaries due to security reasons. All binary distribution listed above are compiled by selected maintainers, who are in charge for all binaries of their platform, respectively.
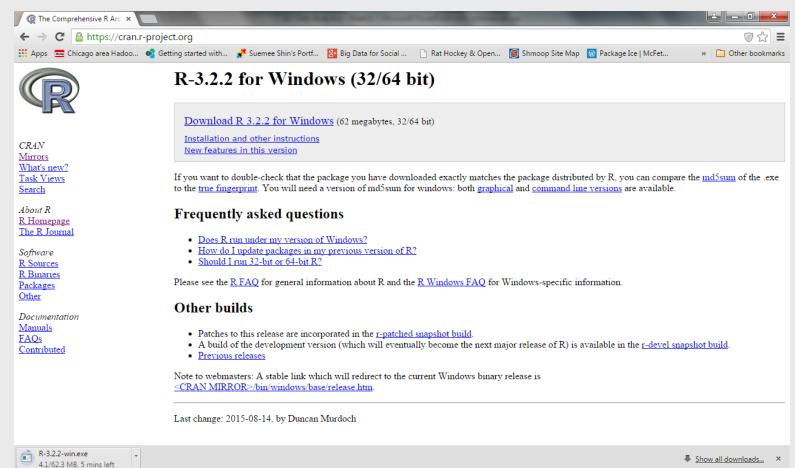
ITM - 527

# Week 11 Topic:
# Getting started with R

◆ R is mainly used when the data analysis task requires standalone computing or analysis on individual servers. It's great for exploratory work, and it's handy for almost any type of data analysis because of the huge number of packages and readily usable tests that often provide you with the necessary tools to get up and running quickly. R can even be part of a big data solution.

◆ When getting started with R, a good first step is to install the amazing RSTUDIO IDE (http://www.rstudio.com/products/rstudio/) . Once this is done, we recommend you to have a look at the following popular packages:

◆ dplyr, plyr and data.table to easily manipulate packages,

◆ stringr to manipulate strings,

◆ zoo to work with regular and irregular time series,

◆ ggvis, lattice, and ggplot2 to visualize data, and

◆ caret for machine learning

**40**

# Week 11 Topic:
# R Pros and Cons

◆ **_Pro: A picture says more than a thousands words:_** Visualized data can often be understood more efficiently and effectively than the raw numbers alone. R and visualization are a perfect match. Some must-see visualization packages are ggplot2, ggvis, googleVis and rCharts.

◆ **_Pro: R ecosystem:_** R has a rich ecosystem of cutting-edge packages and active community. Packages are available at CRAN, BioConductor and Github. You can search through all R packages at Rdocumentation (http://www.rdocumentation.org/).

◆ **_Pro: R lingua franca of data science:_** R is developed by statisticians for statisticians. They can communicate ideas and concepts through R code and packages, you don't necessarily need a computer science background to get started.  Furthermore, it is increasingly adopted outside of academia.

◆ **_Pro/Con: R is slow:_** R was developed to make the life of statisticians easier, not the life of your computer. Although R can be experienced as slow due to poorly written code, there are multiple packages to improve R's performance: pqR, renjin and FastR, Riposte and many more.

◆ **_Con: R has a steep learning curve:_** R's learning curve is non-trivial, especially if you come from a GUI for your statistical analysis. Even finding packages can be time consuming if you're not familiar with it.

**41**

ITM - 527

# Week 11 Topic:
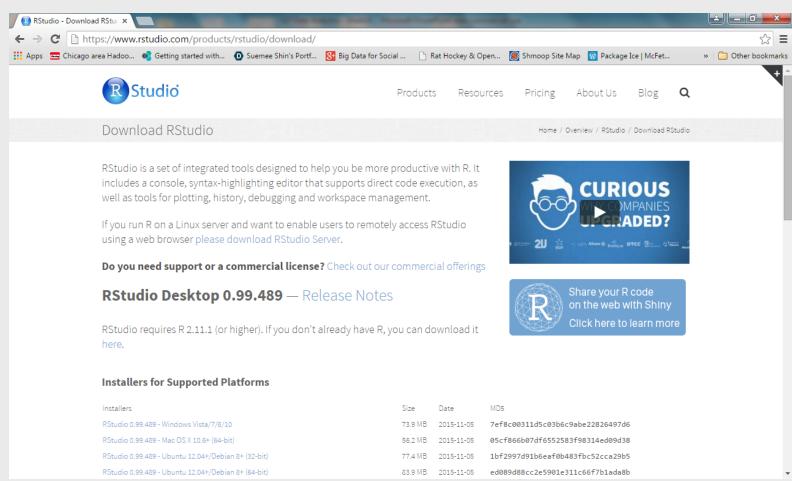# Installing R

http://cran.r-project.org/

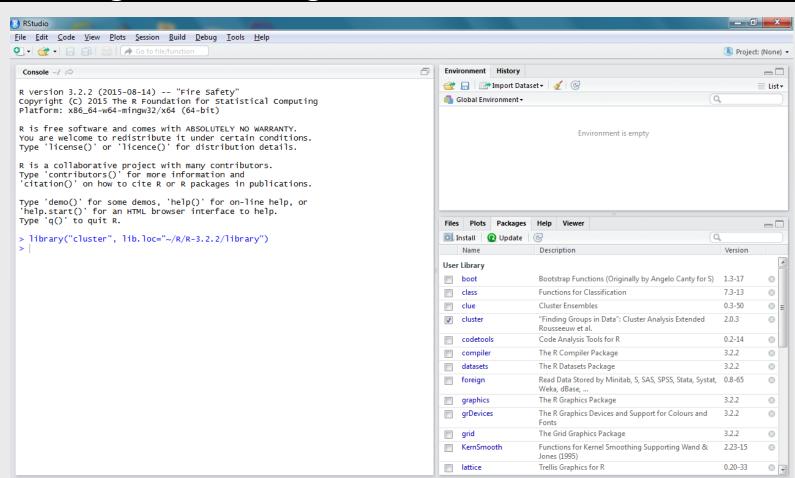# Week 11 Topic:
# (Optional) R Training & References

◆    Revolution Analytics offers good video training on R:
http://www.revolutionanalytics.com/academyr-training-education

◆    An Introduction to R (In Blackboard)

◆    Big Data (In Blackboard)

ITM – 527
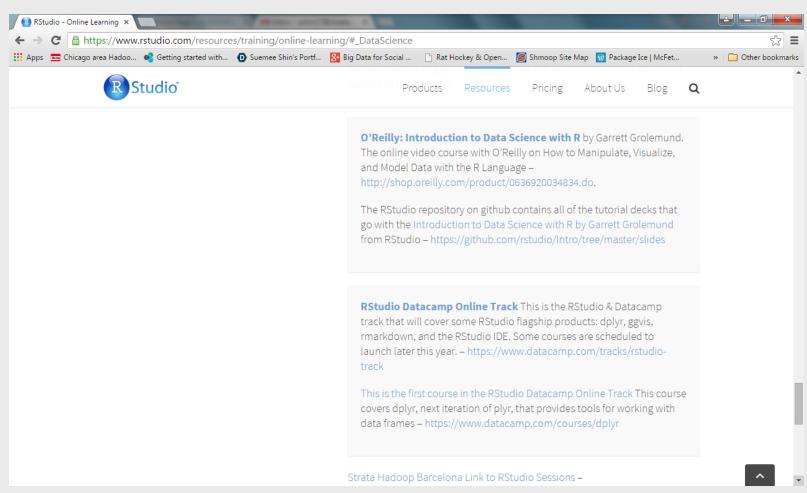
43

ITM - 527

# Week 11 Topic:
# Installing RStudio

# Week 11 Topic:
# Using R Packages

# Week 11 Topic:
# RStudio Resources

ITM - 527

# Week 11 Topic:
# TryR