

ICEAA Workshop 2016

Integrating Cost Estimating and Data Science Methods in R

Josh Wilson and Laura Barker

Atlanta, GA
June 2016

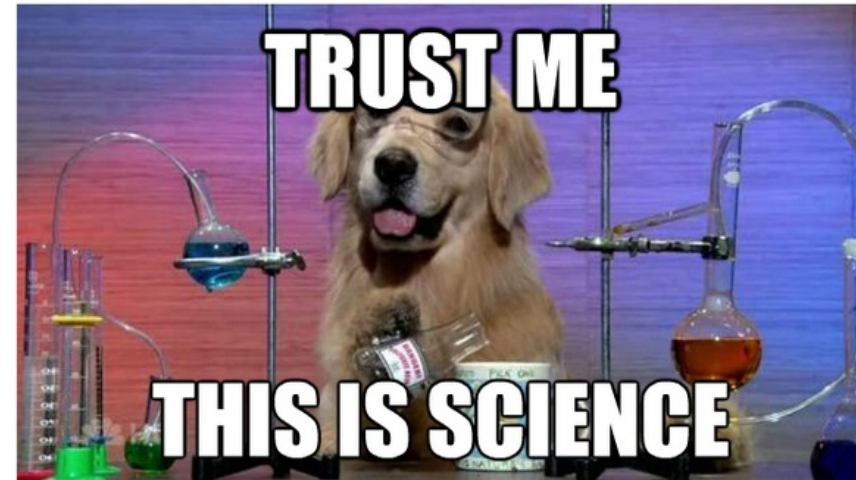
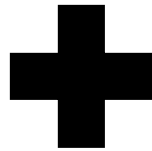
Outline

- 4 Purpose / Overview
- 4 Definition of Terms
- 4 Cost Estimating Challenges
- 4 Example – Installation Cost Analysis in Excel and R
- 4 Benefits (Pros)
- 4 Limitations (Cons)
- 4 Conclusions
- 4 Way Forward

Purpose / Overview

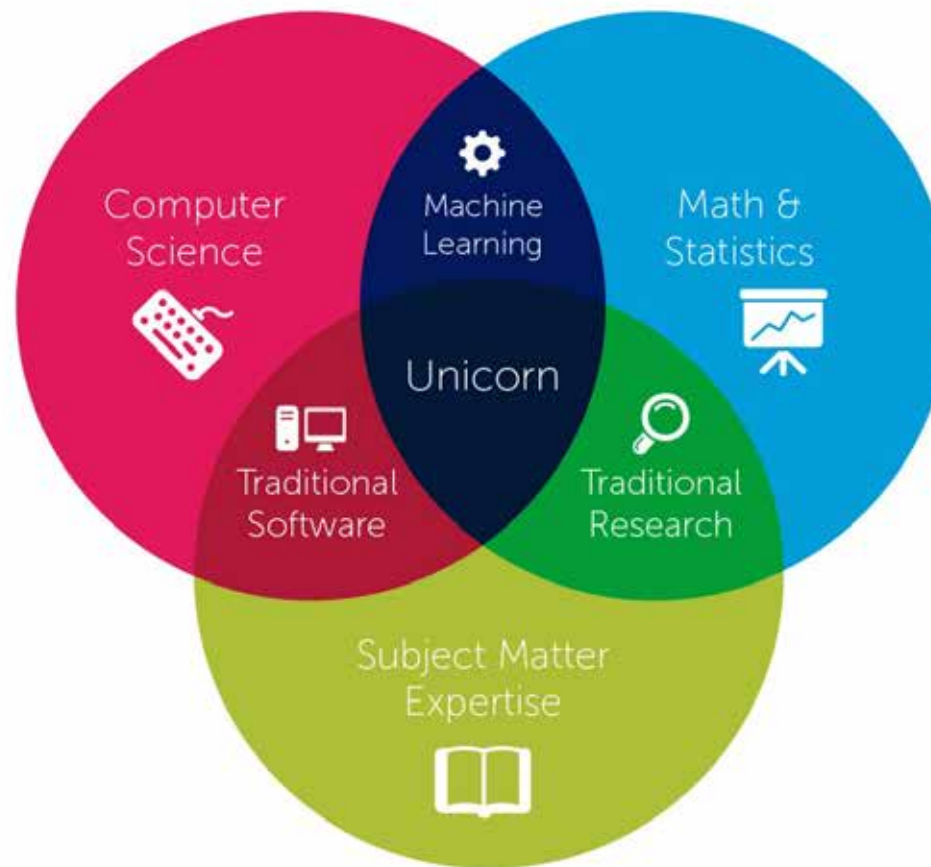
- 4 Data science tools can help address some basic, persistent cost estimating challenges
 - Documentation
 - Traceability
 - Reproducibility
 - Reusability
- 4 Particularly applicable to data collection and analysis stages of the cost estimating process
- 4 Usage of data science tools also opens the door to more advanced data science techniques

What is Data Science?



No, seriously...

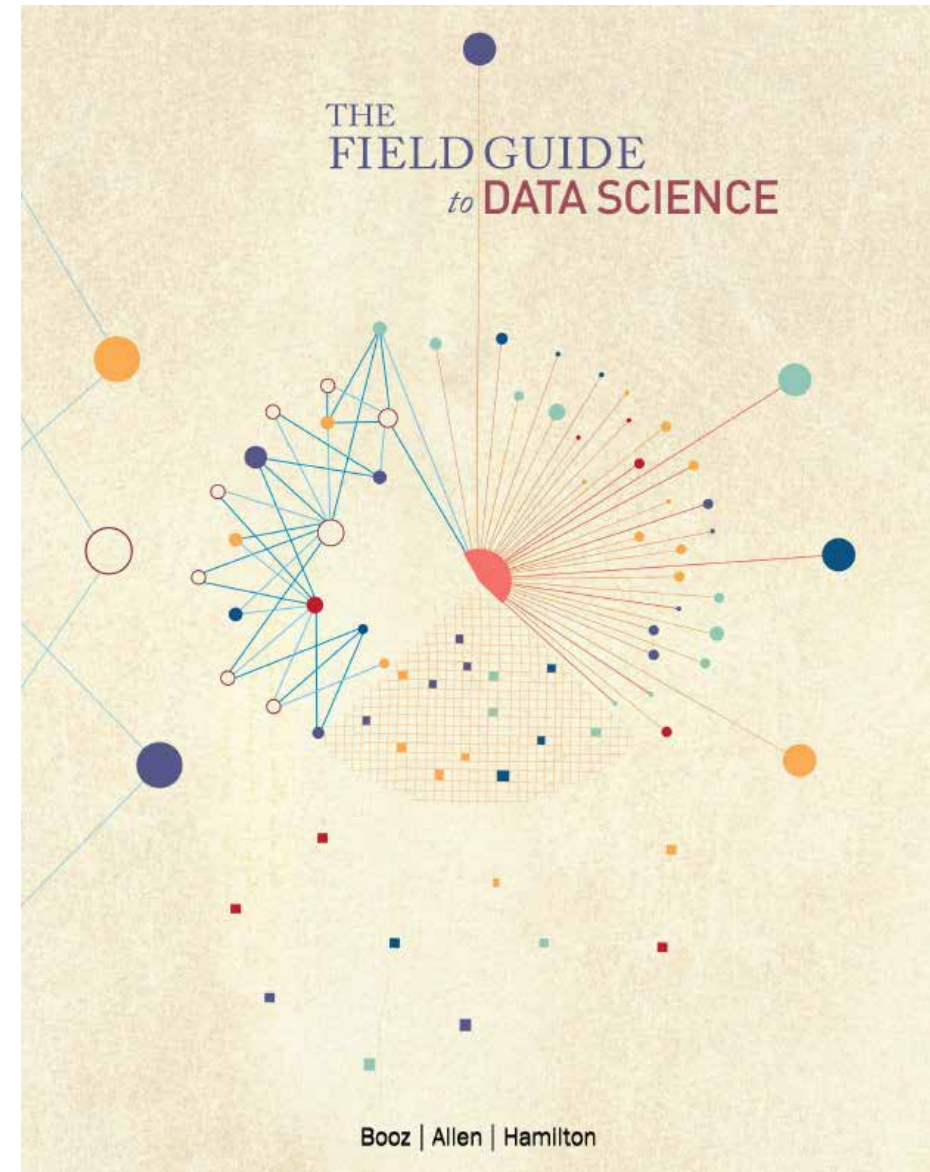
Data Science



Copyright © 2014 by Steven Geringer Raleigh, NC.
Permission is granted to use, distribute, or modify this image,
provided that this copyright notice remains intact.

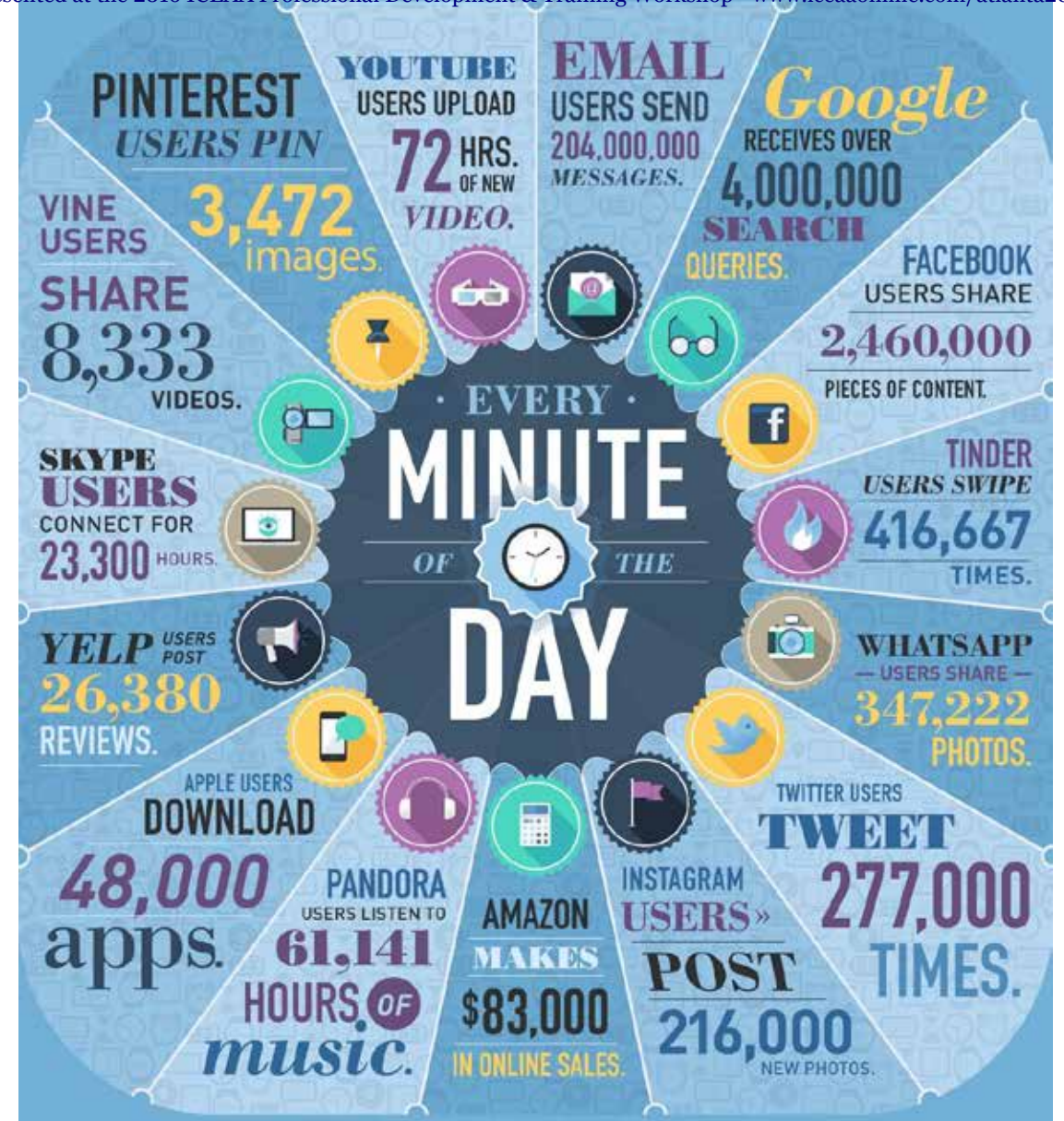
What *is* Data Science?

- 4 “Describing Data Science is like trying to describe a sunset – it should be easy, but somehow capturing the words is impossible... Data Science is the art of **turning data into actions.**” – The Field Guide to Data Science
- 4 “The key word in ‘Data Science’ is not Data, it is Science.” – Jeff Leek
- 4 “I think data-scientist is a sexed up term for a statistician... Statistics is a branch of science.” – Nate Silver



Why do you keep hearing about Data Science?

- 4 POP QUIZ: 90% of the world's data has been created in the past ____ years.
- 4 Data infrastructure advances
 - Everything is connected (Internet of Things)
- 4 Commoditization of computing power
 - Data collection is getting easier and cheaper



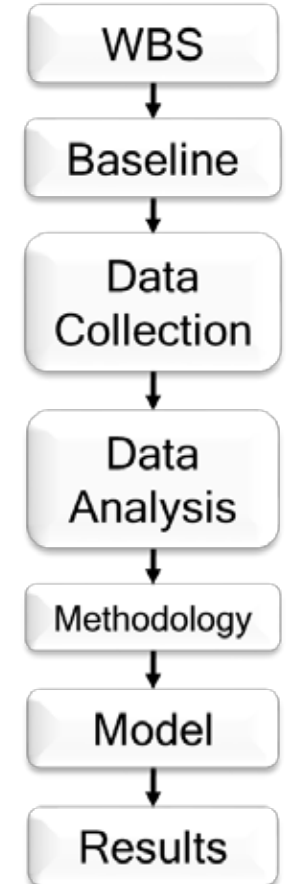
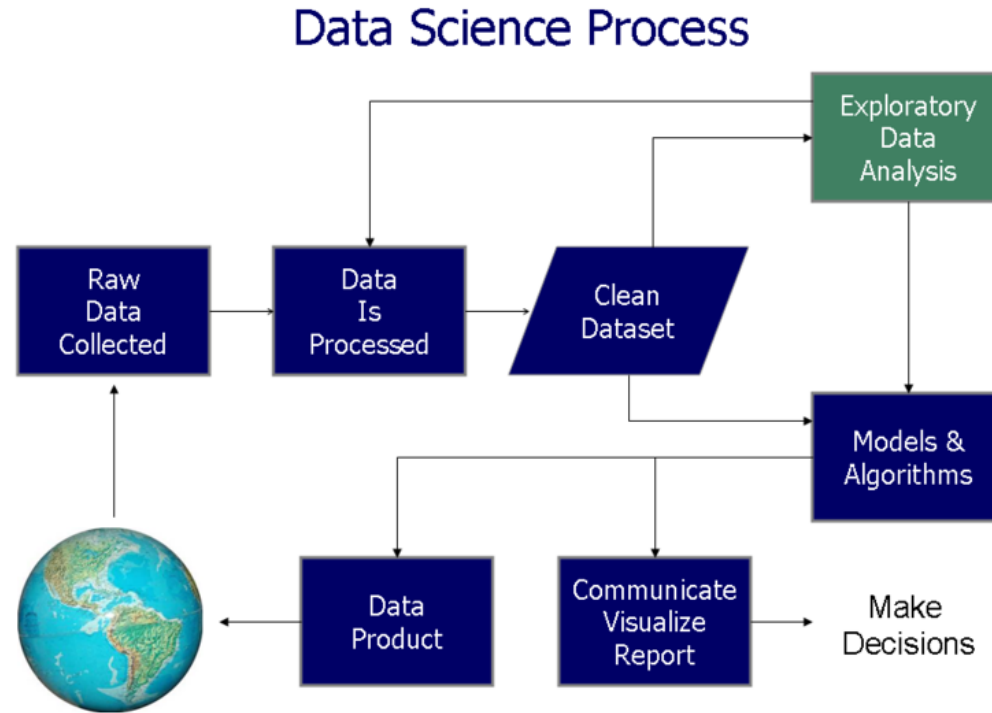
Why are *WE* talking about Data Science?

4 Data science is about gaining insight from data

4 Data science emphasizes:

- Documentation
- Traceability
- Reproducibility
- Reusability
- Statistical rigor

4 Data science process aligns with parts of the cost estimating process



Cost Estimating Challenges

4 Documentation

- Set of explicit instructions to define data sources and analysis steps
- Data collection and analysis process is inherently difficult to document

4 Traceability

- Ability to follow and understand the process used to perform an analysis
- Limited by poor documentation

4 Reproducibility

- Ability to perform an analysis again and get the same result
- Limited by lack of documentation and traceability

4 Reusability

- Ability to apply a completed analysis to a new situation without having to redo everything

How can Data Science enhance Cost Estimating?

Documentation – Set of instructions to define data sources and analysis steps

4 Traditional cost estimating processes and methods:

- 4 Documentation is an extra step to the analysis process
- 4 Various levels of quality
 - Source data and analysis may be included within a cost model (as an Exhibit or Appendix)
 - Source data and analysis may be in a workbook external to main cost model
 - External analysis model could have been lost
- 4 Any required Excel add-ins or macros are not typically identified

4 Cost estimating supplemented by Data Science:

- 4 Analysis scripts (code) are self-documenting
- 4 Imported packages are clearly identified
- 4 Additional documentation features are available
 - R Markdown
 - R Pubs
 - Slidify
 - Python Sphinx

How can Data Science enhance Cost Estimating?

Traceability – Ability to follow and understand the analysis process

4 Traditional cost estimating processes and methods:

- 4 Heavily dependent on having a clean model structure and good model documentation
- 4 Never ending string of Excel precedents / dependents
- 4 Lots of details are easy to miss
 - Custom filters
 - Ad hoc data cleaning
- 4 Excel Stats package output is hardcoded

4 Cost estimating supplemented by Data Science:

- 4 Scripted analysis steps are clearly defined in a sequential manner
- 4 Many popular data science tools are open-source
 - Python
 - R

How can Data Science enhance Cost Estimating?

Reproducibility – Ability to perform an analysis again and get the same result

4 Traditional cost estimating processes and methods:

4 Ability to reproduce analysis dependent on quality of model documentation and traceability

4 Requires manual reproduction of analysis steps

4 Access to Excel add-ins may cause issues

4 Cost estimating supplemented by Data Science:

4 Simply requires re-running the script

4 Required packages imported and loaded automatically

How can Data Science enhance Cost Estimating?

Reusability – Ability to easily apply a completed analysis to a new situation

4 Traditional cost estimating processes and methods:

4 Reuse of analysis in Excel requires reproduction of effort

- Re-application of filters
- Re-navigation of stats package GUIs
- Potential manual changing of formula references

4 Visuals usually need to be created based on specific data

4 Cost estimating supplemented by Data Science:

4 Once created, a scripted analysis is readily applied to new data and/or inputs

4 Less "rework" required to reuse analysis script

- Change source data
- Change inputs

4 Automatically generates visuals

EXAMPLE

Question: “Is the cost to install a system different on various ship classes?”

Example: Analysis of Install Cost Data in Excel

1. Obtain raw data pull
2. Apply filters
 - “Job Type” equals “INSTALL”
 - “Percent” equals “100”
 - “Product” contains text “SYSTEM NAME”
3. Add column to extract ship class info from hull code
 - Use formula “LEFT(*hull code*, FIND(“ “, *hull code*))”
4. Manually copy and paste values to another sheet
 - Manipulate data so that data for each ship class is in a different column
 - Necessary to simplify generation of graphs

The screenshot illustrates the process of filtering and data manipulation in Excel. It shows a data table with columns J through M. The 'Job Type' column (M) is filtered to show 'INSTALL' and 'SHORE'. The 'Percent' column (O) is filtered to show '100'. A 'Custom AutoFilter' dialog box is open, showing the criteria 'contains SYSTEM NAME'. A formula bar shows the formula `=LEFT(K1683,FIND(" ",K1683))` being applied to the 'Hull/Cd' column (K). The resulting data is shown in the table below:

Ship/Cd	Hull/Cd	Ship Cl.	BF/REG
SS MILIU	DDG 0069	" ,K1683))	
SS ENTEF	CVN 0065	CVN	
SS DONA	DDG 0075	DDG	
SS BAINE	DDG 0096	DDG	
SS NITZE	DDG 0094	DDG	THEODOR
SS HUE C	CG 0066	CG	

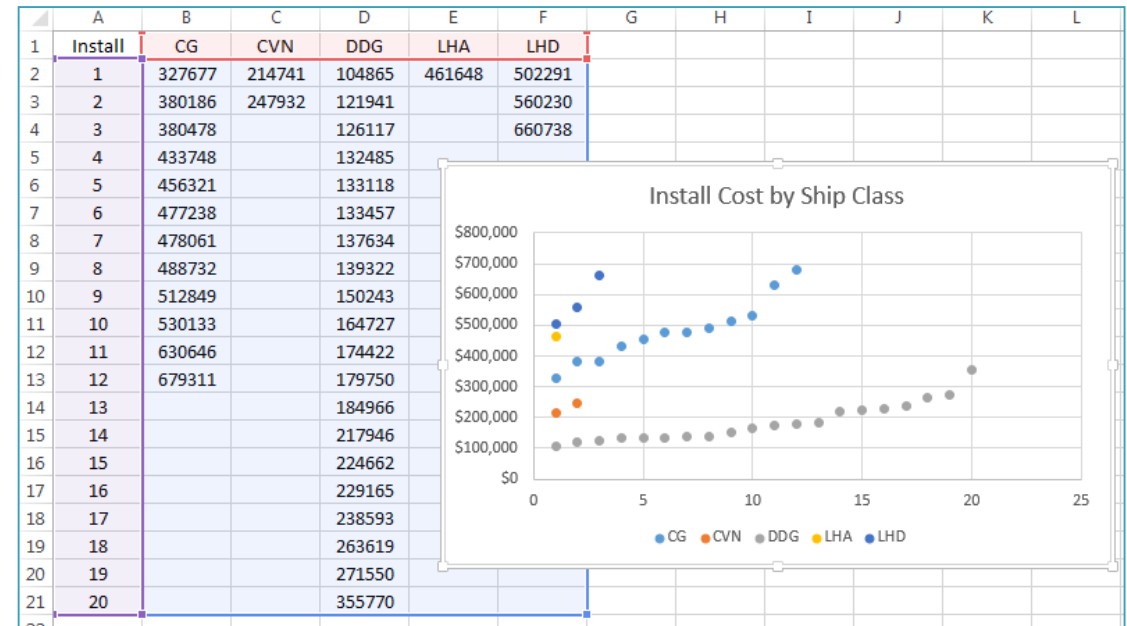
Example: Analysis of Install Cost Data in Excel (cont'd)

5. Scatter plot install cost

- Separate ship class columns put each ship class in an individual data series
- Note that some GUI navigation steps here were skipped / i.e. not documented

6. Formulate hypotheses based on visual inspection of data

- Install costs for DDG and CVN appear indistinguishable
- Install costs for CG, LHA, LHD also appear indistinguishable
- Install costs for {CG, LHA, LHD} > {DDG, CVN}

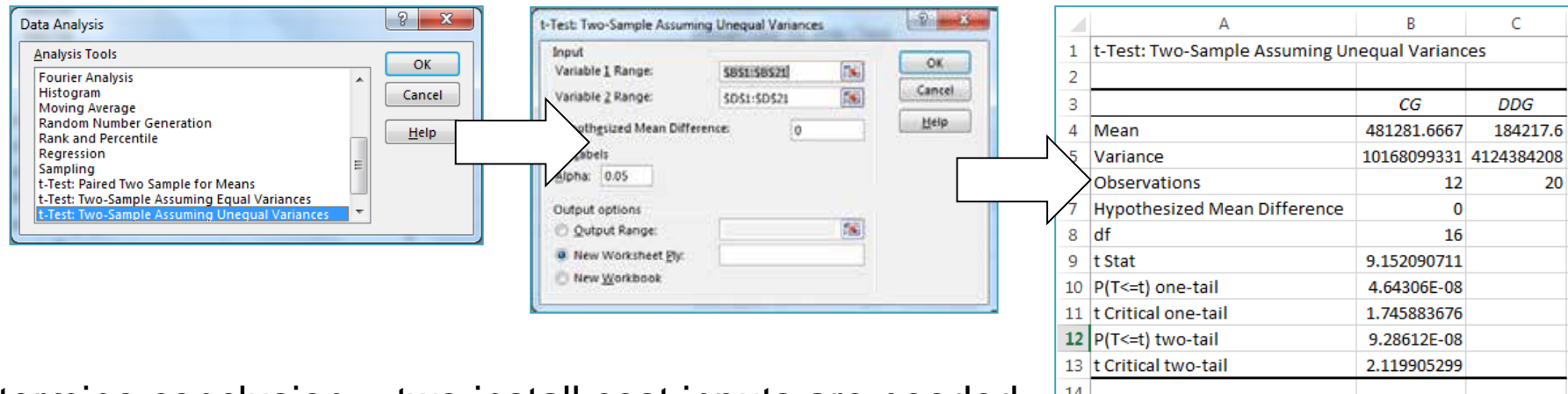


* DDG = Destroyer, CVN = Carrier, CG = Cruiser, LHA / LHD = Types of Amphibious Assault Ships

Example: Analysis of Install Cost Data in Excel (cont'd)

7. Hypothesis testing via multiple t-tests

- Need $(N*(N+1))/2$ individual pairwise tests to compare N ship classes (5 ship classes = 10 comparisons)
- Below is t-test of CG vs. DDG cost



8. Determine conclusion – two install cost inputs are needed

- CVN / DDG
- CG / LHA / LHD

Example: Analysis of Install Cost Data in Excel (cont'd)

4 Questions the cost analyst may be asked during presentation of results

- What is the supporting data?
- What were the steps of the analysis?
- Does the analysis process make sense?
- Can the analysis results be reproduced?
- Can the analysis results be easily updated when new data is available?
- Can the analysis be applied to other systems?

- Documentation
- Traceability
- Reproducibility
- Reusability

Example: Analysis of Install Cost Data in Excel (cont'd)

4 Documentation issues

- Lots of small, manual steps are easy to overlook
- GUI navigation steps require screenshots
- Required add-ins and macros not typically identified

4 Traceability issues

- Copying and pasting specific values loses any reference to source data
- Hardcoded t-test results eliminate visibility to test input data

4 Reproducibility issues

- Difficult due to documentation and traceability issues

4 Reusability issues

- Requires re-filtering, re-copying and pasting values, re-navigation of stats GUIs
- Manual changes to update plots

Example: Analysis of Install Cost Data in R

4 What is R?

- Free, open source programming language
- Used for data manipulation, data analysis, statistical modeling, data visualization, and predictive modeling
- Users include data scientists, statisticians, analysts
- Has gained substantial popularity in recent years
- Easily extended via packages

4 Analysis in R is a set of instructions

- Script to read, manipulate, and analyze data

```

1 > ##### INSTALL / LOAD PACKAGES #####
2 if(!require("dplyr"))
3   install.packages("dplyr")
4
5 if(!require("ggplot2"))
6   install.packages("ggplot2")
7
8 > ##### READ DATA #####
9
10 ## change wd to the folder in which you have "SPIDER_san2.csv" saved
11 ## ("SPIDER_san2.csv" is the source data file)
12
13 wd <- "C:/Users/534300/Desktop/work Files/Data Science/DSWG/Install Ana
14 setwd(wd)
15 all_data <- read.csv("SPIDER_san2.csv") %>% tbl_df
16
17 > ##### FILTER DATA #####
18
19 ## filter inputs
20 filter_job_type <- "INSTALL"
21 filter_percent <- 100
22
23 ## product names have been sanitized by replacing with NES game titles
24 filter_product <- "Mega Man"
25 # filter_product <- "The Addams Family"
26 # filter_product <- "The Incredible Crash Dummies"
27 # filter_product <- "M.U.L.E."
28 # filter_product <- "Kung-Fu"
29
30 ## filter data
31 filtered_data <- all_data %>%
32   filter(job.type == filter_job_type) %>%
33   filter(Percent == filter_percent) %>%
34   filter(grepl(filter_product, Alt.Brief.Product))
35
36 ## add ship class from hull code
37 hull_code <- filtered_data$Hull.Code
38 hull_code2 <- strsplit(as.character(hull_code), " ")
39 hull_code3 <- sapply(hull_code2, function(x) x[1])
40 filtered_data$Ship.Class <- hull_code3
41 remove(hull_code, hull_code2, hull_code3)
42
43 > ##### EXPLORATORY PLOTS #####
44
45 > ##### HYPOTHESIS TESTING #####

```

Example: Analysis of Install Cost Data in R (cont'd)

4 Documentation

- Source data clearly identified
- Imported packages clearly identified
- All steps clearly documented

4 Traceability

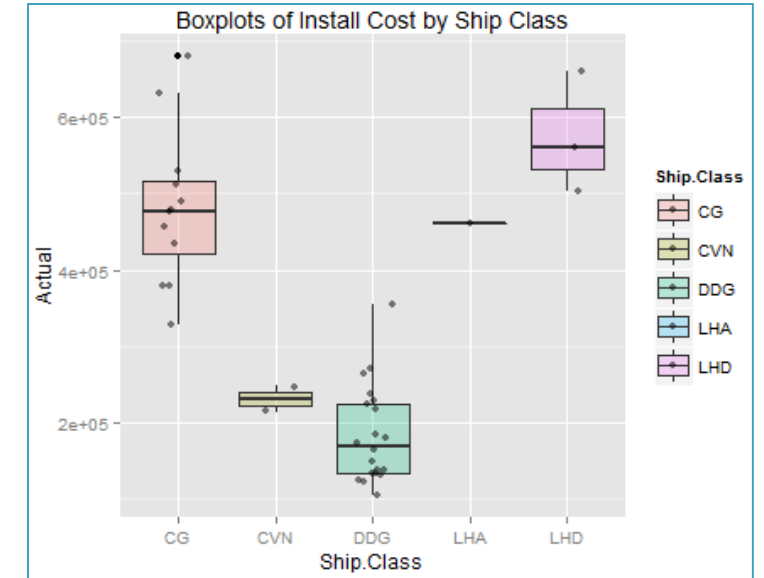
- Less than 100 lines of code
- R is open-source (as is Python)

4 Reproducibility

- Re-running script on source data reproduces the analysis

4 Reusability

- Only requires changing source data and/or filter parameters
- Automates generation of plots and t-test results



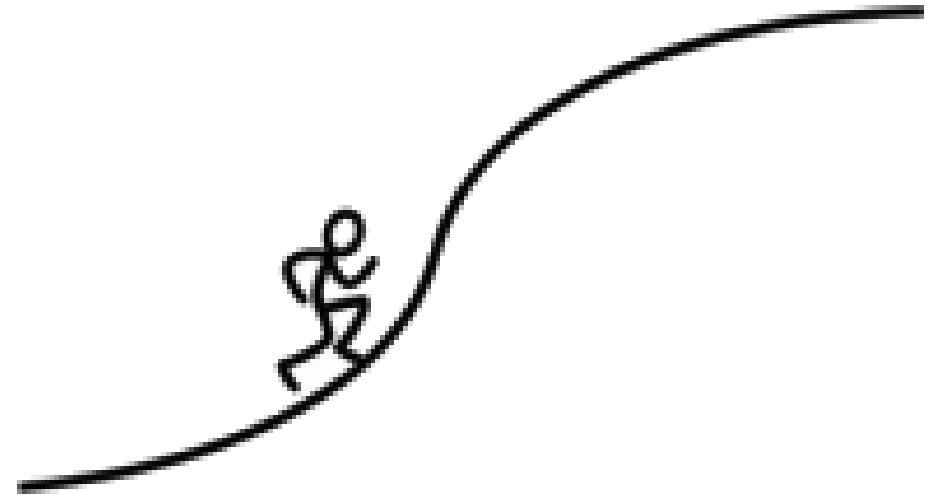
	Ship_Class.1	Ship_Class.2	p_value
1	DDG	CG	7.394639e-08
2	DDG	LHD	8.599988e-03
3	CVN	CG	4.013784e-05
4	CVN	LHD	1.109978e-02

Benefits

- 4 Code is self-documenting
- 4 Additional documentation features are available (R markdown, R pubs, Slidify, Sphinx)
- 4 Some of the most popular data science tools are open source (R, Python)
- 4 Function code in open source languages can be inspected
- 4 Traceability is maximized, making data analysis process more transparent
- 4 Reproducibility is automatic once the analysis script is complete
- 4 Reusing an analysis requires less effort
- 4 More time can be spent improving existing analysis scripts in lieu of performing same analysis on new data
- 4 Robust statistical functionality is built-in
- 4 Opens the door to more advanced analysis techniques

Limitations

- 4 Learning curve
- 4 Unfamiliarity with tools may impact ability to communicate analysis process
- 4 Commonality of Excel versus data science tools
- 4 Every cost estimate is unique
- 4 Garbage in, garbage out



Conclusions & Way Forward

- 4 Data science can be applied anywhere there is data, including cost estimating
- 4 Data science tools and practices can be used to address some common cost estimating challenges
- 4 Data science tools will not replace cost estimating tools or methods, but are more appropriate in some cases
- 4 Utilizing data science tools opens the door to advanced analysis techniques
 - Clustering algorithms to identify truly analogous programs or efforts
 - Textual analysis of contractor monthly status reports to predict cost overruns
 - k-NN (k Nearest Neighbors) to address missing info in messy data sources
 - Interactive and dynamic visuals (R Shiny package)

Questions?

Josh Wilson
Associate

Booz | Allen | Hamilton

Booz Allen Hamilton Inc.
1615 Murray Canyon Road
Suite 900
San Diego, CA 92108
Tel (619) 278-3855
Mobile (619) 820-6226
wilson_joshua@bah.com

Laura Barker
Senior Consultant

Booz | Allen | Hamilton

Booz Allen Hamilton Inc.
1615 Murray Canyon Road
Suite 900
San Diego, CA 92108
Tel (619) 680-4857
Mobile (770) 990-6905
barker_laura@bah.com

BACKUP

What is R?

- 4 R is a free, open source programming language and software environment for statistical computing and graphics
- 4 The R language is widely used among statisticians and data miners for developing statistical software and data analysis, and has gained substantial popularity in recent years

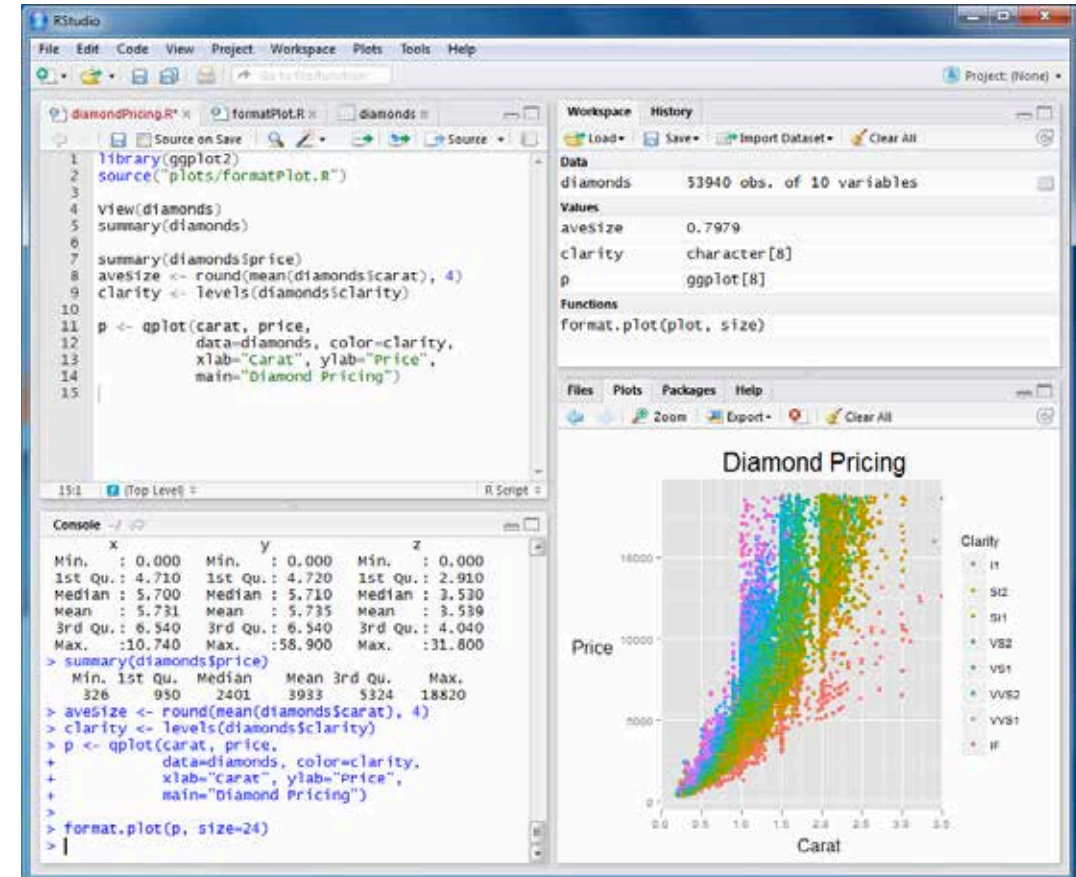


What is RStudio?

4 RStudio is a free Integrated Development Environment (IDE) for R that has added functionality and more user-friendly interface

4 Includes:

- console
- syntax-highlighting editor that supports direct code execution
- easier file navigation
- tools for plotting, history, debugging, and workspace management



Why choose R over other Data Science tools?

4 Free

4 Open source

4 Large user base

4 (Relatively) easy to learn

4 Built-in tools and packages for statistical analysis

- Linear and nonlinear modelling
- Classical statistical tests
- Time-series analysis
- Classification

4 Data visualization capabilities