

01-Introduction

Practical Machine Learning (with R)

UC Berkeley

Spring 2016

Topics

- Introduction
- Tools and Environment
- Exercise
- Introductions (continued)
- Data Science, Machine Learning and Opportunities
- Elite Coding
- Exercises
- Machine Learning



INTRODUCTIONS



Me (Personally)



Shameless Plug

My Skills

- R /Python Programmer (>15 years)
- Machine Learning (>15 years)
- DevOps
- Researcher and Writer : Machine Learning, Clinical Medicine, Chemistry, Finance

Education

- UC Berkeley → (UT Austin)→ UC Santa Barbara → UC Berkeley
- Post-graduate: UC Berkeley, Stanford

Professional Experience

- Lawrence Berkeley National Lab, Allianz, Open Data
- Sept. 2010 Founded Decision Patterns

Professional Interests

- Machine Learning / Statistics
- High Performance Computing
- Applied Statistics and Visualization
- Management of Data Organizations



(Decision Patterns)



Shameless Plug

Decision Patterns

- Founded 2010
- Bring together complementary skills for data strategy:

Acquisition * Organization * Storage Access * Utilization

- Our Model
 - Service Consulting
 - Not a start-up -- no VC funding
 - Use consulting margins from to niche products
- Our Customers
 - Financial Services, Retail, Entertainment, Food, Communications, Defense, Environmental Sciences



What do I like *most* about what I do?





BEST
THING

We get to work on a

- variety of problems,
- with a variety of technologies
- in a variety of fields



What do I like *least* about what I do?



WORST
THING

We have to work on a

- variety of problems,
- with a variety of technologies
- in a variety of fields



TOOLS AND ENVIRONMENT



EXERCISE: SET-UP TOOLS AND ENVIRONMENT

- ⇒ Install **R** → **CRAN**
- ⇒ Install **R Studio Desktop™** (IDE)
- ⇒ Install **git**

- ⇒ Create **github** account
 - Send name, student id to christopher.brown@berkeley.edu



GIT



Git / GITHUB / Source Tree Workflow

What is it?

A source control tool
(and **process**) to promote
collaborative
development

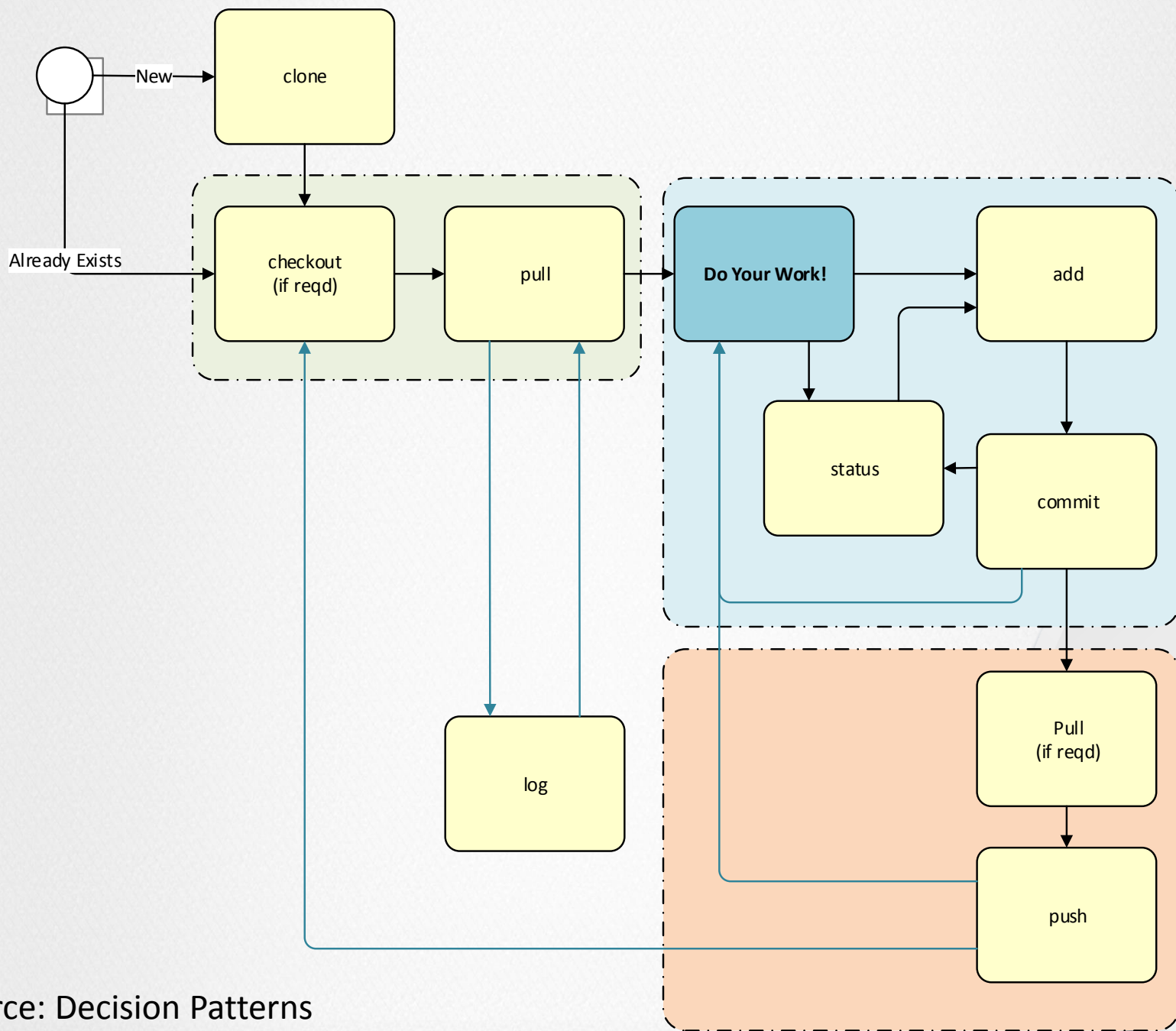
Features

- Distributed
- Each clone contains complete history
- Ability to return to
- Branch in and merging

Interfaces

- Github
- Source Tree™
- *R Studio*
- *command-line*





source: Decision Patterns

GIT COMMANDS

- **Repo(sitory)**: location where files are stored. If different from original source: “fork”
- **Branch**: Copy of code that can be independently worked on.
- **checkout**: Change to specific branch/commit.
- **add**: Tell which files to “stage” (accept) commit. Done on a per-file basis.
- **commit**: accept changes.
- **pull**: Retrieve changes from remote repository
- **push**: Send committed change to remote repository
- **log**: review history of commits
- **status**: review “staged” status

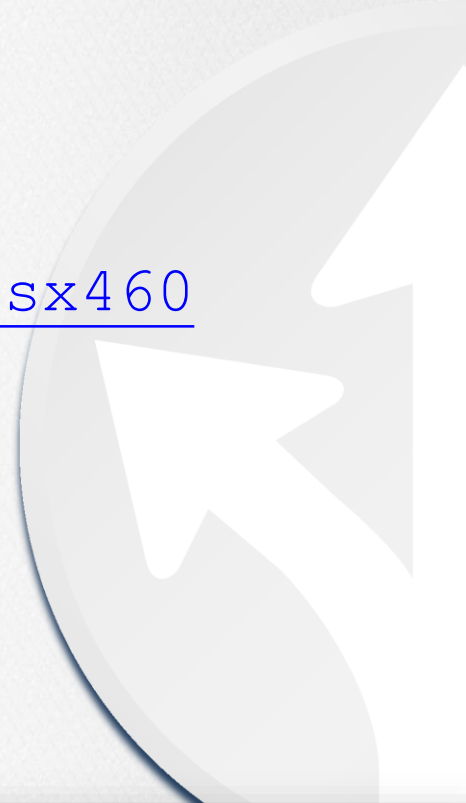


EXERCISES

- ➔ Create github account
- ➔ Send account log-in, student id to christopher.brown@berkeley.edu

- ➔ Clone class repository

```
git clone https://github.com/csx460
```



You?

DISCUSSION OF INDIVIDUAL GOALS ?



You

- ➔ How many of you are students? Professionals?
- ➔ How many have
 - > 1 year using R?
 - > 3 years?
 - > 5 years?
- ➔ How many use R as your principal data.science tool?
- ➔ How many use
 - Python
 - Julia
 - SAS or SPSS
 - Spark/Scala
 - Java
- ➔ Ever spend too much time debating which technology fits?

Class / Objectives

Theory

- Distinguish fundamental aspects of machine learning algorithms
- Build (train) machine learning models
- Evaluate (score) machine learning models
-

Practice

- Frame problems to make the suitable for solution via machine learning
- Collaborate in a group using tools for collaborative/social programming
- Generate high quality, graphical and textual results
- Deploy machine learning models to operations

Required Text

Applied Predictive Modeling

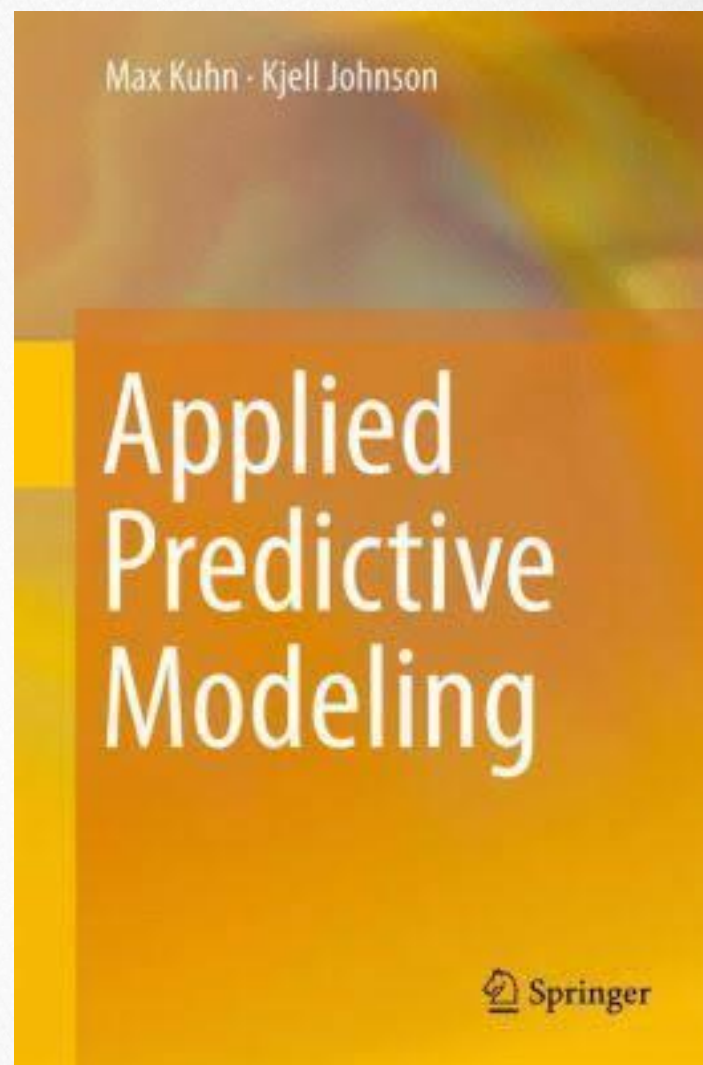
ISBN-13: 978-1461468486

ISBN-10: 1461468485

Kuhn, Max and Johnson, Kjell

Springer Science+Business

2013



Additional Resources

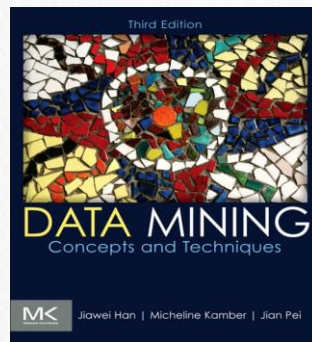
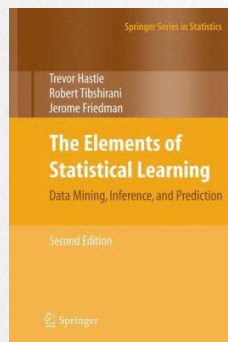
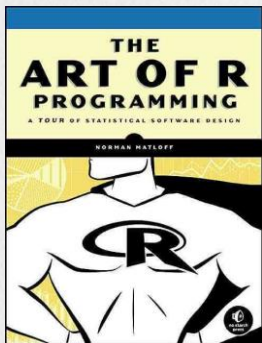
Texts

(not used in this class)

The Art of R Programming
by Norm Matloff

Elements of Statistical Learning
by Hastie, Friedman, Tibshirani

Data Mining Concepts and Techniques
by Han, Kamber, Pei



Online

- ➔ CRAN
 - [Packages](#)
 - [Task Views](#)
- ➔ [Metacran](#) (r-pkg.org)
- ➔ [Stackoverflow.com](#)
- ➔ [r-bloggers.com](#)
- ➔ [Advanced R Programming](#)
- ➔ [Github](#)



CONTACTS / COORDINATES

- ➔ Christopher Brown
christopher.brown@berkeley.edu
checked once / day (mornings)
phone #
- ➔ Class Website
 - <https://github.com/CSX460>
- ➔ Google Group: CSX460



GRADING

⇒ ~8 Weekly Exercises (80%)

- Exercises are **Rmarkdown** in the github
- Due at the beginning of class each week
- Submitted via **github** commits
 - Please email me your **github** login
 - ! Github commits are timestamped
- Answers reviewed in class
- Work on them in class, time-permitting

⇒ Class Participation (20%)

⇒ Attendance is Mandatory

- no unexcused absences.



**** PARTICIPATE ****



RMARKDOWN



RMARKDOWN

What is it?

- Simple text mark-up syntax
 - that supports the markdown standard
 - And allows incorporation of R analysis and graphical output
- Are assignments will be done in Markdown ...
 - Simply put your answers in the space provided
 - → Demonstration



ML OVERVIEW

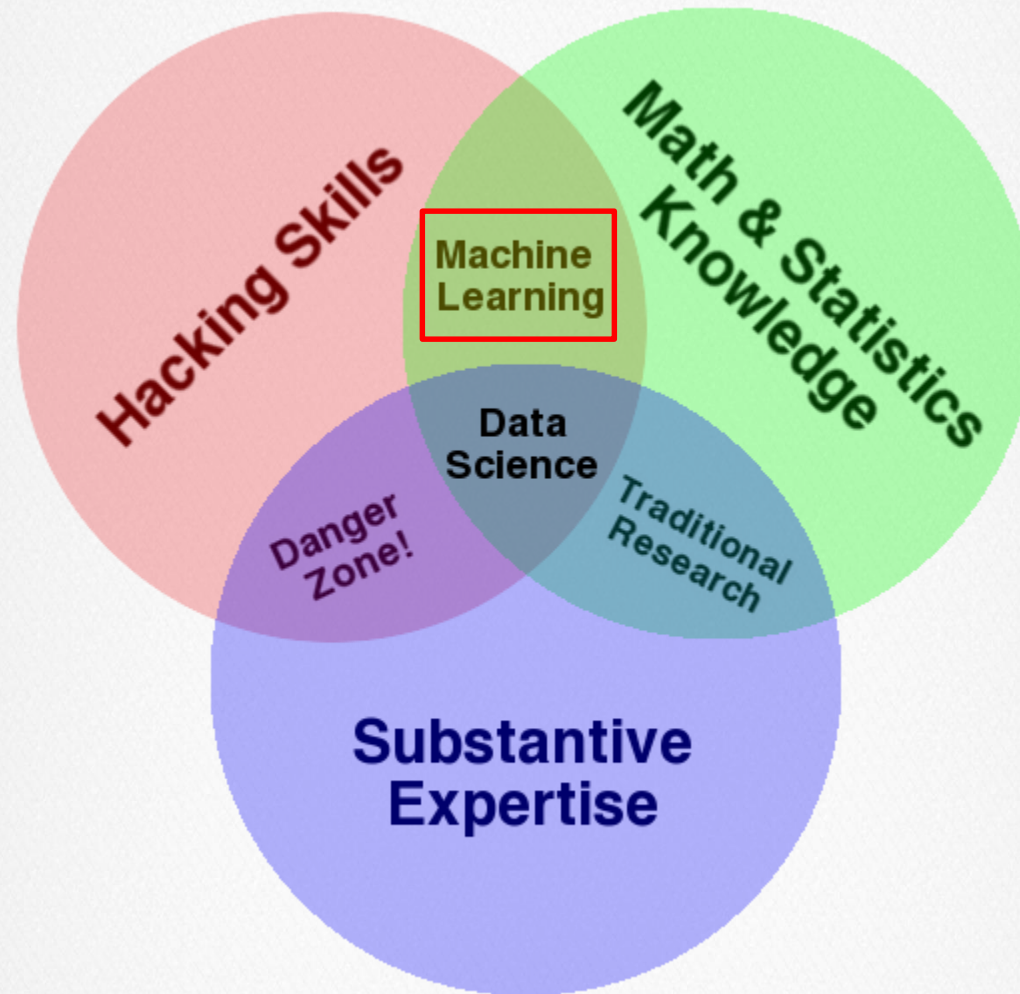


EXAMPLE OF ML ALGORITHM(S)

- Spam Filter
- handwriting recognition (svm)
- Traffic engineering (lights)
- Weather prediction
- Sentiment analysis (social media)
- Netflix Recommender
- Fraud detection (Visa)
- Imaging processing
- Intrusion detection
- Self-driving cars



Data Science Venn Diagram



Ref. <http://drewconway.com/zia/2013/3/26/the-data-science-venn-diagram>

MACHINE INTELLIGENCE 2.0

AGENTS

PROFESSIONAL	PERSONAL	OS INTERFACES
 DigitalGenius 	 	 Google Now

AUTONOMOUS SYSTEMS

AIR	GROUND	SEA	INDUSTRIAL

ENTERPRISE

SECURITY / FRAUD	HR / RECRUITING	SALES	MARKETING	CUSTOMER SUPPORT	INTERNAL INTEL	MARKET INTEL

PLATFORMS

RESEARCH / AGI	FULL STACK	MACHINE LEARNING	INDUSTRIAL IOT	AUDIO	VISION	DATA ENRICHMENT

INDUSTRIES

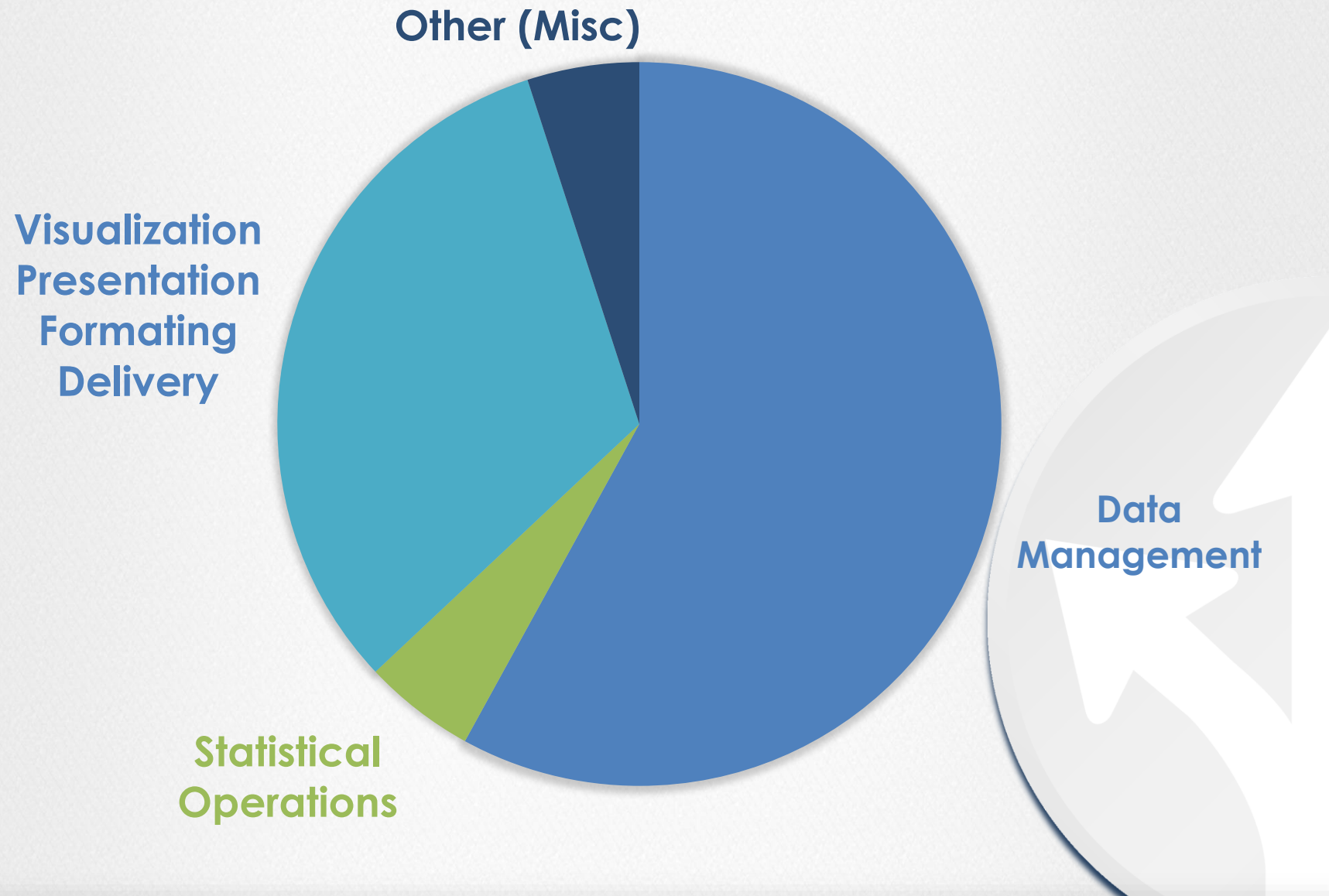
ADTECH	AGRICULTURE	FOR GOOD	RETAIL FINANCE	LEGAL	MATERIALS & MFG	HEALTHCARE

INDUSTRIES (CONT'D)

EDUCATION	TRANSPORT & LOGISTICS	INVESTMENT FINANCE	DATA SCIENCE	MACHINE LEARNING	OPEN SOURCE

SHIVONZILIS.COM/MACHINEINTELLIGENCE

BREAKDOWN OF CODE TASKS



ELITE CODING



ELITE CODING / 1

→ Follow Established Design Patterns

CREATIVITY IS GENERALLY A BAD THING

Goal	Description	R Packages
Ad hoc analysis	Create a process	ProjectTemplate, Rmarkdown, knitr
Package Development	Create a package	Rstudio, Roxygen2, devtools
Application : Interactive	Web application	Shiny, OpenCPU Javascript
Application : Automated	Code to be scheduled or called as an event	Rscript (R -e), optigrab, crontab

ELITE CODING / 2

→ Adopt standards: python™

- Hadley Wickham's **Advanced R** style guide
<http://adv-r.had.co.nz/Style.html>
- Decision Patterns Style Guide
- Do **NOT** follow Google's coding convention
- Cf. Python's Standards
 - PEP-8 Naming and Formatting
 - PEP-257 Documentation
 - PEP-20 Readability

→ Use version control:



- Github, Bitbucket, Gitlab
- Best GUI: Atlassian Sourcetree

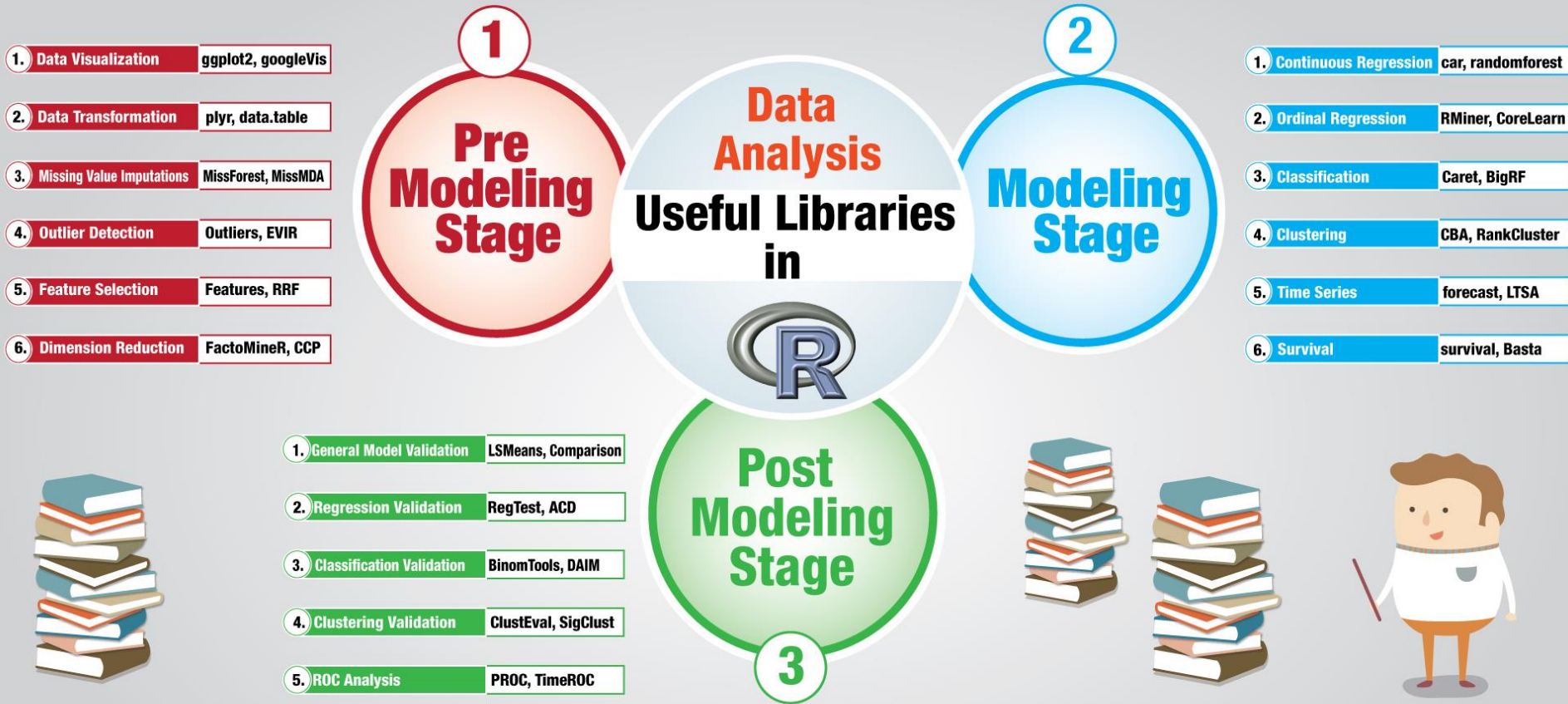
Commit early and often.

→ Use Agile Methods

- Track issues: JIRA, Github, Gitlab

Good PM is worth every penny.

```
>install.packages("package name")
```



Other Libraries

A. Improve performance: Rcpp, parallel

B. Work with web: XML, jsonlite, httr

C. Report results: shiny, RMarkdown

D. Text Mining: tm, twitterR

E. Database: sqldf, RODB, RMongo

F. Miscellaneous: swirl, reshape2, qcc

USEFUL R PACKAGES

```
> install.packages ("package-name")
```

- ⇒ **ML Framework:** *caret* (Classification and Regression Training)
- ⇒ **Pipe operators:** *magrittr* (*pipeR*, *backpipe*) (*shiny*)
- ⇒ **Tables:** *data.tables*, *dplyr*
- ⇒ **Visualization:** *ggvis*, *ggplot2*
- ⇒ **Reporting:** *knitr*, *rmarkdown*, *shiny*



EXERCISES IN CLASS



QUESTION 1

What is machine learning?

A formal **process** for building a **model**



QUESTION 2

What is a model?

a ***function*** that ***estimates*** a ***response*** associated with (a set of) known ***predictors***

$$\hat{y} = f(\vec{x})$$



QUESTION 3: WHAT ARE THE PROPERTIES OF f

- ⇒ Should be easy* to evaluate
- ⇒ Takes a one or more values of inputs
- ⇒ Yields a single output value for each input
- ⇒ Output, \hat{y} , should be “close to” observed values, y :

$$\hat{y} \sim y$$

* Computational cheap/efficient

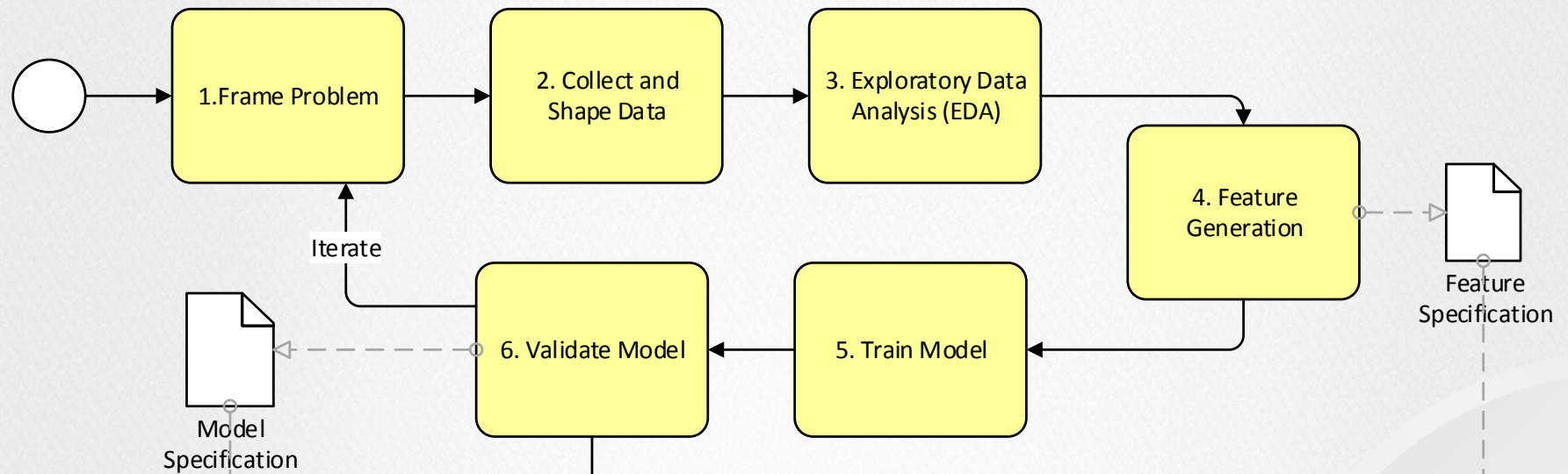


QUESTION 4

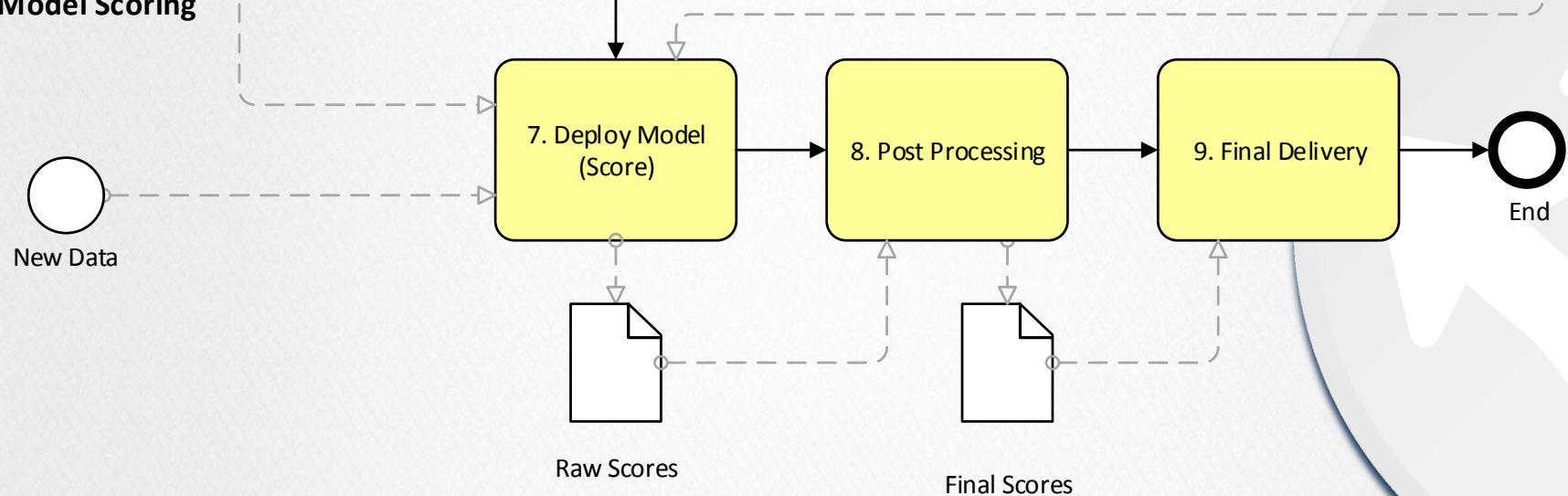
How do we find f ?



Model Training



Model Scoring



APPENDIX



CLASS OVERVIEW : 1

- Introduction to R, setting up the ML developers environment
 - Installing R
 - Installing R Studio
 - Installing packages from CRAN, Bioconductor and Github
 - Exercises



CLASS OVERVIEW : 2

- ➔ Fundamentals of Machine Learning
 - Machine learning overview
 - Regression and classification
 - Supervised, unsupervised, and semi-supervised
 - Algorithm types and requirements
 - Exercises



CLASS OVERVIEW : 3

- ➔ Linear Regression (2 sessions)
 - OLS Regression
 - Data partitioning
 - Model evaluation and tuning
 - Exercises



CLASS OVERVIEW : 4

- ⇒ Logistic Regression
 - Logistic Regression
 - Exercises



CLASS OVERVIEW : 5

- Advanced Techniques: Partitioning Methods
 - CART/Regression Trees
 - Clustering
 - K Nearest Neighbors
 - Exercises



CLASS OVERVIEW : 6

- Advanced Techniques: Partitioning Methods
 - CART/Regression Trees
 - Clustering
 - K Nearest Neighbors
 - Exercises



CLASS OVERVIEW : 7

➤ Advanced Techniques

- Bagging
- Bagged Trees / Random Forests
- Exercises



CLASS OVERVIEW : 8

- Advanced Techniques: Boosting
 - Boosting
 - Neural Networks
 - Support Vector Machines
 - Exercises



CLASS OVERVIEW : 8

⇒ Deployment

- Diving into the data lake
- Optimization
- Delivery and Production



CLASS OVERVIEW : 9

⇒ Final Lecture

- Exercises
- Exam





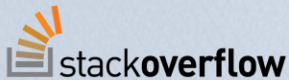

**WHY R?
WHY NOT PYTHON? ... JULIA? ...
SCALA? ...MATLAB?**

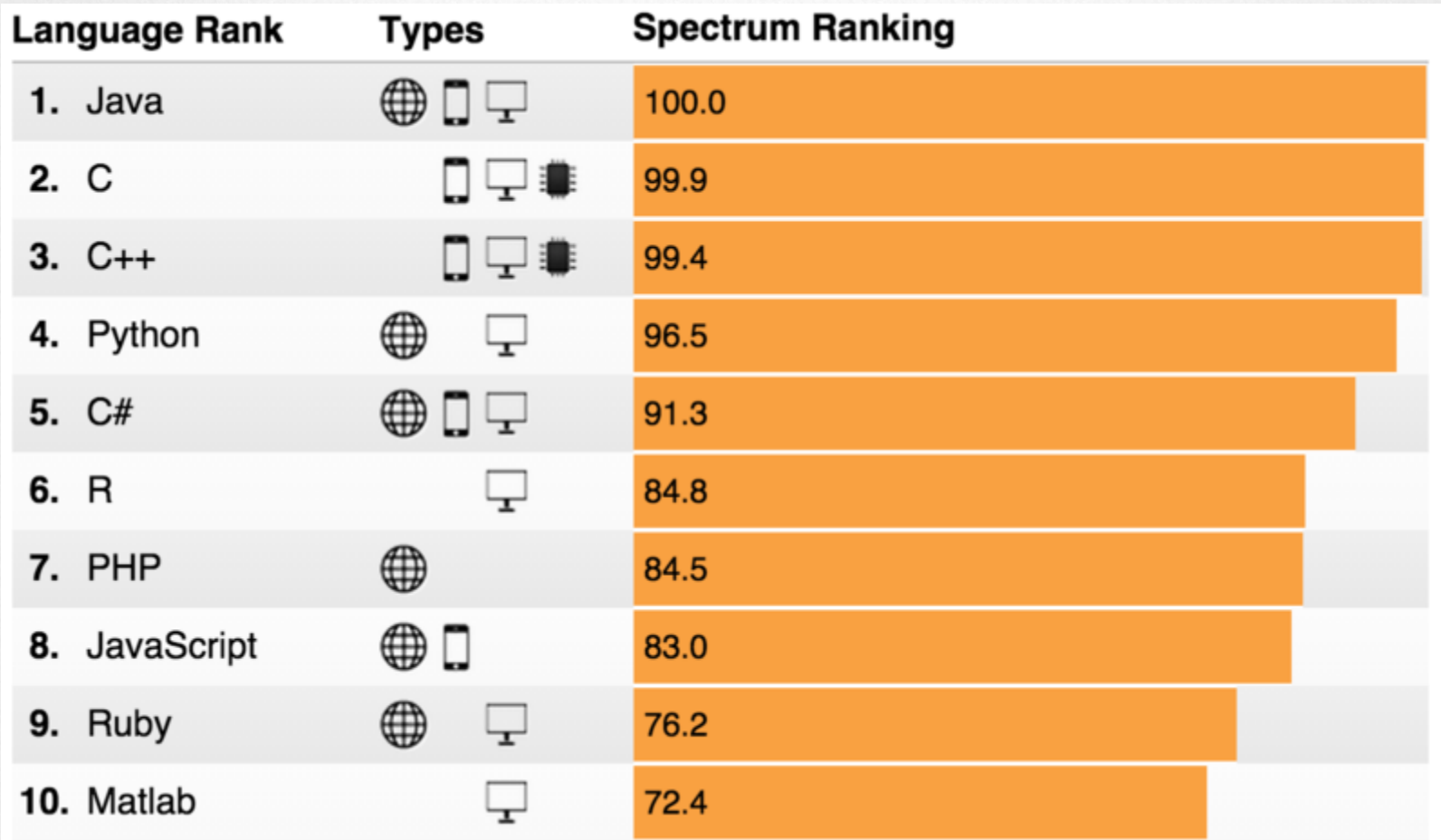


Popularity

2015-06-04



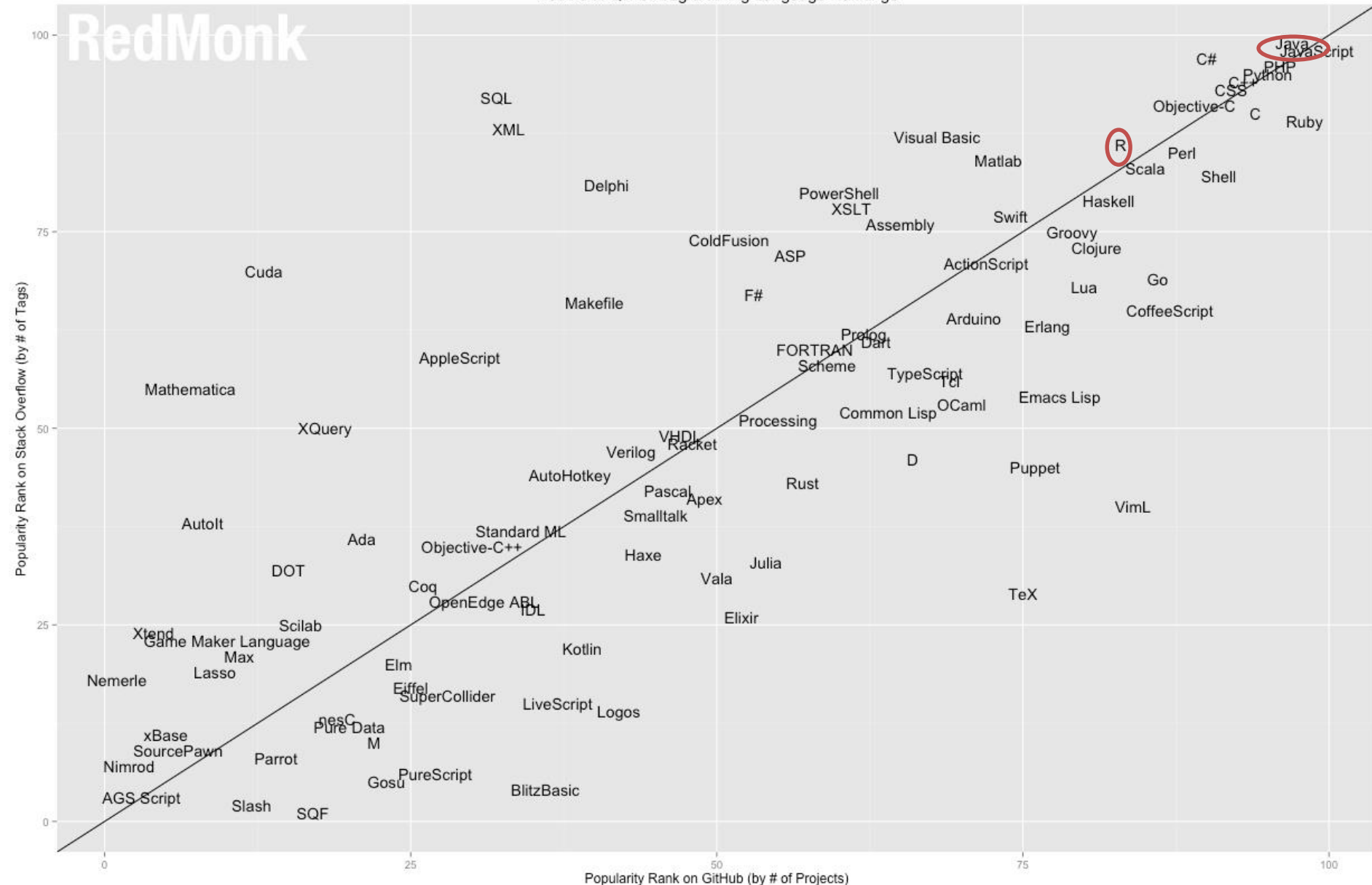
	PyPI	CRAN
Packages	60,806 Packages 35+ updates / day	6,727 package 20+ updates/day
Popularity (Tiobe)	6 th Rank, +0.67% 	12 th Rank, +1.06% 
 stackoverflow	430,604	93,943
 github SOCIAL CODING	549,014	87,306



Ref. <http://spectrum.ieee.org/computing/software/the-2015-top-ten-programming-languages>

RedMonk Q115 Programming Language Rankings

RedMonk





Opinions

Learning Curve

Easier esp. if coming from OO
background

Steeper.
More, dedicated

Code Maintainability

Better package system,
fewer name clashes

Better documentation
Generally less code req'd

Performance

Higher, extensible through
Cython, C, C++

Rcpp

Code expressiveness

Hack to extend operators
Lazy evaluation

Domain Specific
%x% syntax used widely
Non-standard evaluation

Dedicated Web
Frameworks

Translucent(?)

Shiny

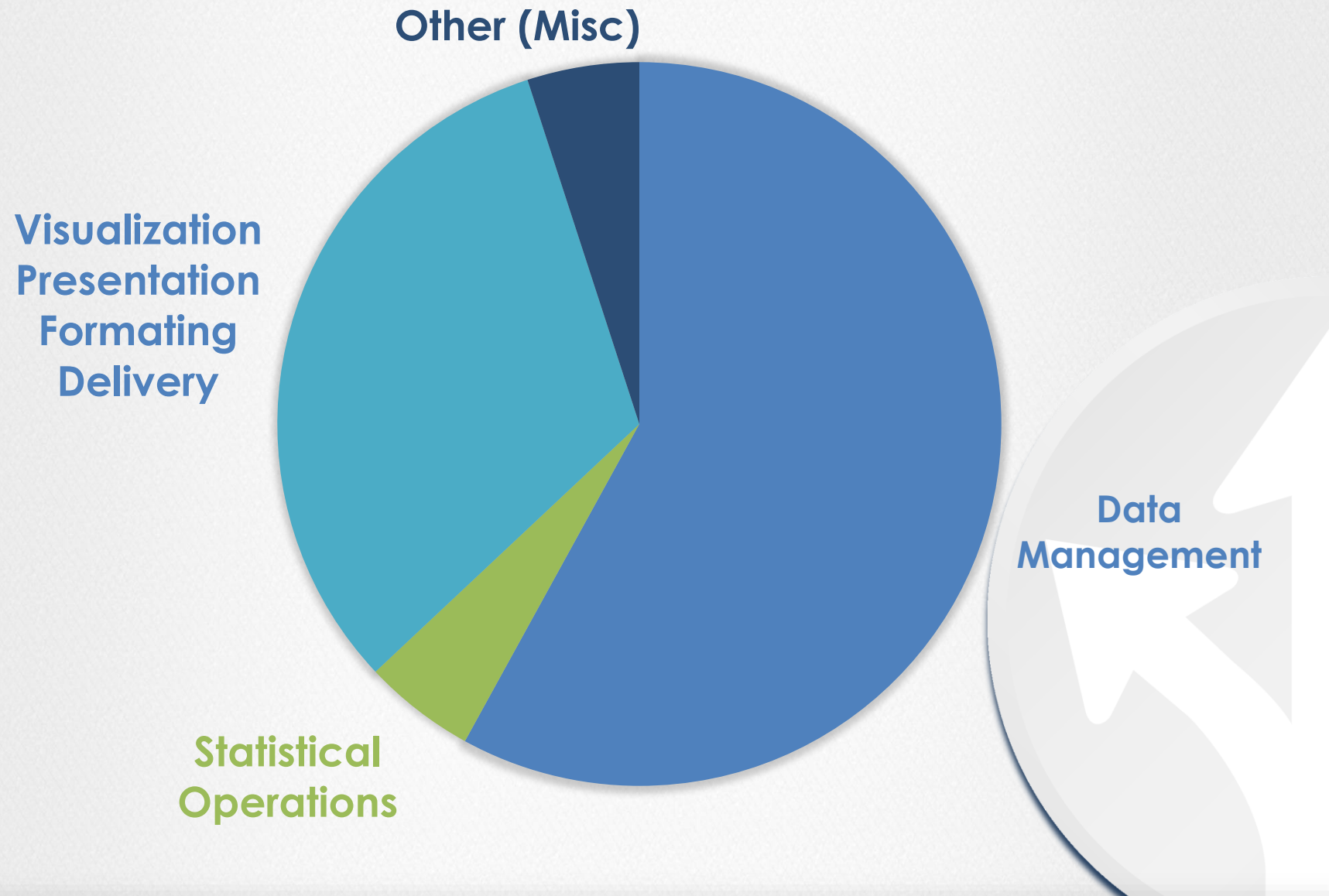
Domain Feature
completeness

Rmarkdown, Reproducible Research,
ProjectTemplate

Vendor Entrenchment

Windows Azure, Oracle, MicroStrategy,
Birst, Tableau, Oracle

BREAKDOWN OF CODE TASKS



R ADVANTAGES

- ⇒ Functional / Vectorized
- ⇒ Dedicated IDE: **Rstudio**
(REPL/Interactive Programming)
- ⇒ **CRAN** and **BioConductor**
- ⇒ **Shiny**
- ⇒ **Domain Specific Abstractions**
 - `data.frame` / `data.table` / `dplyr`
 - model formula
 - `purrr`



R Limitations

- ⇒ Slow
- ⇒ In-memory
- ⇒ Not-scalable



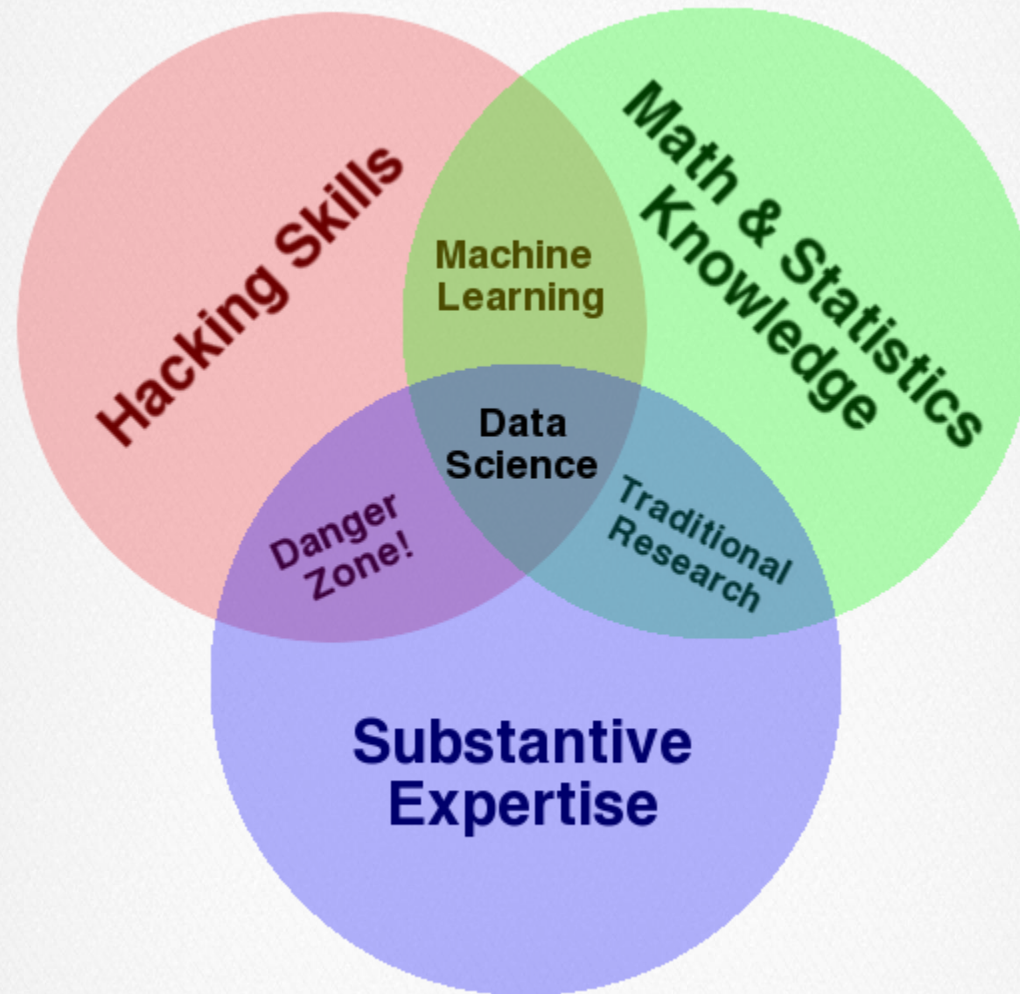
What about ...



DATA SCIENTIST OUTLOOK 2015



Data Science Venn Diagram



Ref. <http://drewconway.com/zia/2013/3/26/the-data-science-venn-diagram>



BusinessIntelligence.com & DOMO
PRESENT



THE WORLD NEEDS DATA SCIENTISTS



IF YOU ARE A MATH- OR DATA-DRIVEN INDIVIDUAL LOOKING FOR THE PERFECT CAREER FIT, look no further than data science. Due to the ongoing explosion of big data, companies have more information at their fingertips than ever—and not enough people who can make sense of it all. This reality has created a big market for quantitative analysts and individuals who can put massive amounts of data into perspective. Take a look.

Source: <http://venturebeat.com/2013/11/11/data-scientists-needed/>

CAREERS IN DEMAND



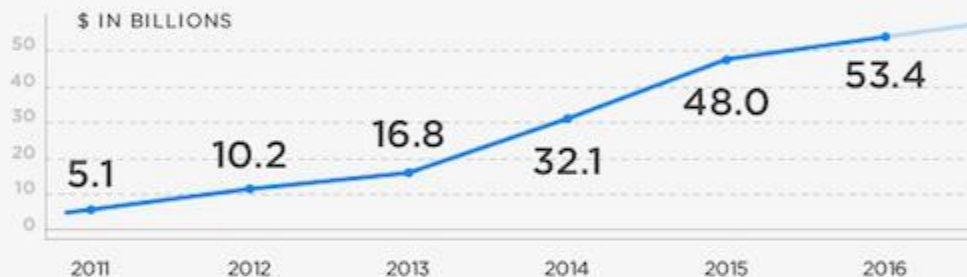


Currently the job market seeks
140,000–190,000
DATA SCIENTISTS TO FILL
OPEN POSITIONS.

IN ADDITION,
1.5 million
data literate managers will need to
be retrained or hired to meet needs.

EXPLAINING THE SUDDEN NEED FOR DATA SCIENTISTS

These scientists don't just happen to be getting far more job offers without reason. Today's modern business needs to manage far more data than ever before, and few have the talent on staff for the job. **Projections indicate that the market will experience meteoric growth in the next several years.**



The Big Data
Market Forecast

Conclusion: With so much activity going on in the big data space and new data touch points being measured every day, there will be an increasing need for data-driven individuals within organizations to make sense of it all. Is that data-savvy person you?



COMPETITION



Much of work will not be done
in traditional worker



H₂O



Google Prediction API

INNOVATION



Spoils go to those who make products
from repeatable processes

The price for analytics is falling ...

