

01-Introduction

Practical Machine Learning (with R)

UC Berkeley

Fall 2015

Topics

- Introductions (Me, You & This Class)
- Data Science Outlook 2015 and Opportunities
- Why R?
- Elite Coding
- Tools and Environment



INTRODUCTIONS



Me (Personally)



Shameless Plug

My Skills

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

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 managing data:

Acquisition * Organization* Storage Access * Utilization

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



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



You?



DISCUSSION OF INDIVIDUAL GOALS ?



Class / Objectives

Theory

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

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

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



GRADING



PROJECT (50%)

- Identify one machine learning problem you want to tackle. It could be from existing UCI data sets, your work, interest, Kaggle ... or I can give you one.
- Frame the problem
- Build Features
- Review linear model and cart
- Build Model
- Deploy



TODAY 2015-09-14

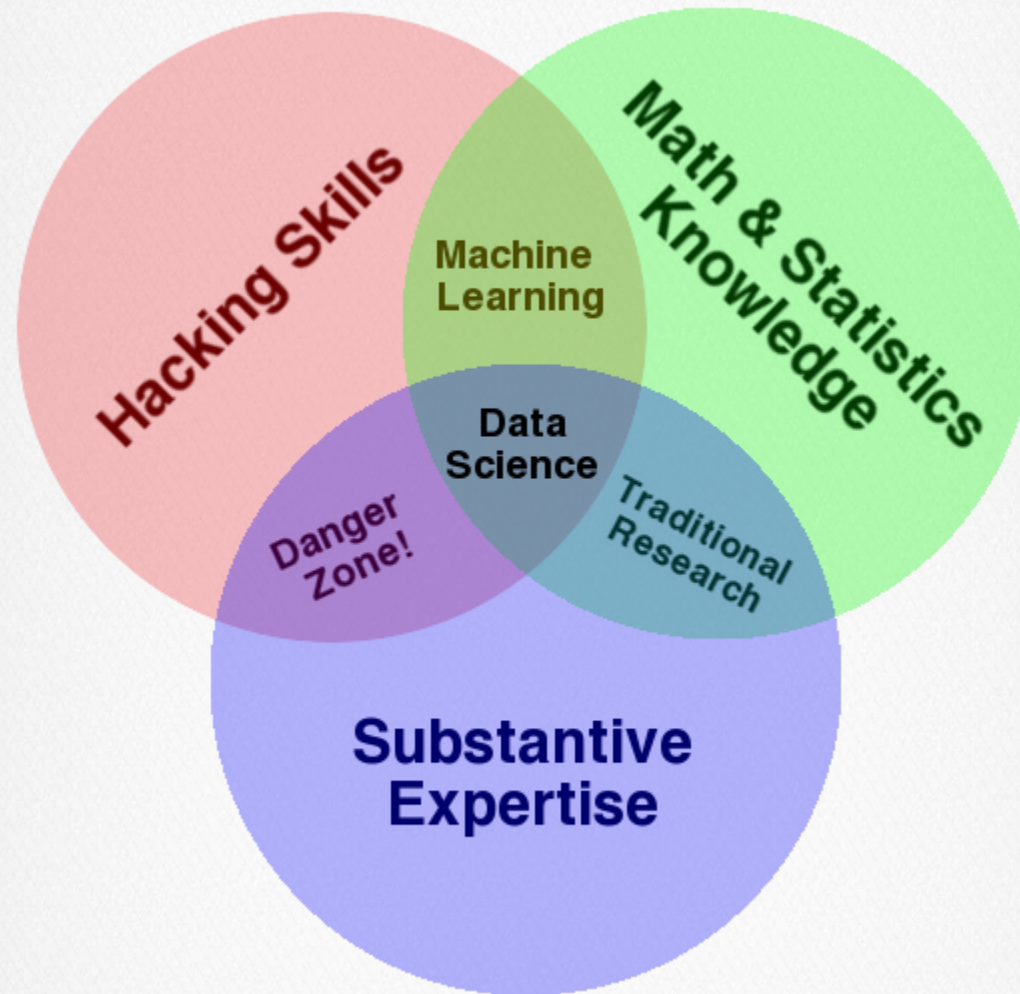
- Data Scientist / Machine Learning Outlook
- R Language
- Effective Programming



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



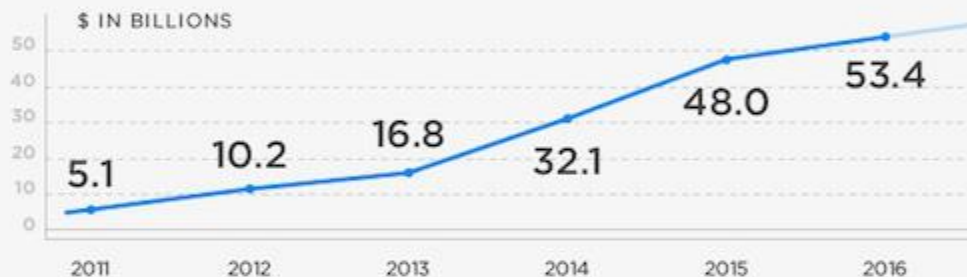


Currently the job market seeks
140,000–190,000
DATA SCIENTISTS TO FILL
OPEN POSITIONS.

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?

Machine Intelligence LANDSCAPE

CORE TECHNOLOGIES

ARTIFICIAL INTELLIGENCE

IBM WATSON MetaMind
Numenta ai-one
Cycorp Research nano
Reactor SCALED INFERENCE

DEEP LEARNING

Vicarious
facebook
Google
SKYMINJ
Baidu
ersatz
SignalSense

MACHINE LEARNING

rapidminer context
data2o G DATA
LiftIgniter Azure ML yhat
GraphLab Alpine AYASDI

NLP PLATFORMS

cortical.io idibon
LUMINOSO wit.ai
Maluba

PREDICTIVE APIS

AlchemyAPI MINDOPS
Google big indico
ALGORITHMIA Expect
PredictionIO Labs

IMAGE RECOGNITION

clarifai MADBITS
DNNresearch DEXTRO
VISENZE lookflow

SPEECH RECOGNITION

GRIDSPACE
popUP archive
NUANCE

RETHINKING ENTERPRISE

SALES

Preact
RelateIQ
CLARABRIDGE
infer
AVISO
NGDATA
FRAMED
ATTENTIV causata

SECURITY / AUTHENTICATION

CROSSMATCH
EYEVERIFY
CYLANCE
conjur
BITSIGHT
bionym

FRAUD DETECTION

sift science
ThreatMetrix
Brighterion
CSOCURE
feedzai
VERAFIN

HR / RECRUITING

TalentBin
predikt
gild
entelo
Connectifier
hiQ
CONEXUS

MARKETING

brightfunnel
CommandIQ
RADIUS
Telapart
bloomreach
AIRPR
people pattern
Freshplum

PERSONAL ASSISTANT

Siri
Cortana
tempo
KASISTO
VIV
Google now
cleversense
Robnlabs
fuse machines
CLARA LABS

INTELLIGENCE TOOLS

ADATAD
Palantir
Quid
Digital Reasoning
FirstRain

RETHINKING INDUSTRIES

ADTECH

METAMARKETS
dstillery
rocketfuel
YieldMo
ADBRAIN

AGRICULTURE

BLUE RIVER
ceresimaging
THE CLIMATE CORPORATION
TerraVision
KONIGSBERG
tule

EDUCATION

Declara
coursera
KNEWTON
kidaptive

FINANCE

Bloomberg
alphasense
Dataminr
FinGenius
KENSHC
minettbrook
BINATIX

LEGAL

Lex Machina
COUNSELYTICS
JUDICATA
DiligenceEngine
brightleaf
RAVEL
Brevia

MANUFACTURING

SIGHT MACHINE
MICROSCAN
IVISYS
BOULDER IMAGING

MEDICAL

Parzival
Genescient
grand round table
transcriptic
ZEPHYR
bina
TUTE

OIL AND GAS

kaggle
biota
TACHYUS
Flutura

MEDIA / CONTENT

Outbrain
newsle
ARRIA
SAILTHRU
wovii
NarrativeScience
Prismatic
ai

CONSUMER FINANCE

affirm
iVenture
Bill GUARD
LendUp
LendingClub
Kabbage

PHILANTHROPIES

DataKind
thorn
DATA GUILD

AUTOMOTIVE

Google
Continental
Cruise
Tesla
Mobileye

DIAGNOSTICS

enlitic
lumiaata
3SCAN
ENTONIS

RETAIL

BAY SENSORS
PRISM SKYLABS
celect
euclid

RETHINKING HUMANS / HCI

AUGMENTED REALITY

usable intelligence
APX
blippar
META
layar

GESTURAL COMPUTING

THALMICLABS
omek
LEAP
eyeSight
GestureTek
nod

ROBOTICS

Intel
iRobot
anxi
LIQUID ROBOTICS
SoftBank
Boston Dynamics
Artificial Robotics

EMOTIONAL RECOGNITION

affectiva
BEYONDERBAL
EMOTIENT
cogito

SUPPORTING TECHNOLOGIES

HARDWARE

NVIDIA
XILINX
QUALCOMM
NERVANA
TeraDeep
rgetti

DATA PREP

TRIFACTA
tamr
Paxata
Alation

DATA COLLECTION

diffbot
kimono
CrowdFlower
Cinnotate
WorkFusion
import io



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 ...





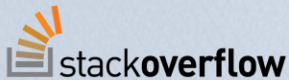

WHY R?
WHY NOT PYTHON? ... JULIA? ...
SCALA? ...MATLAB?

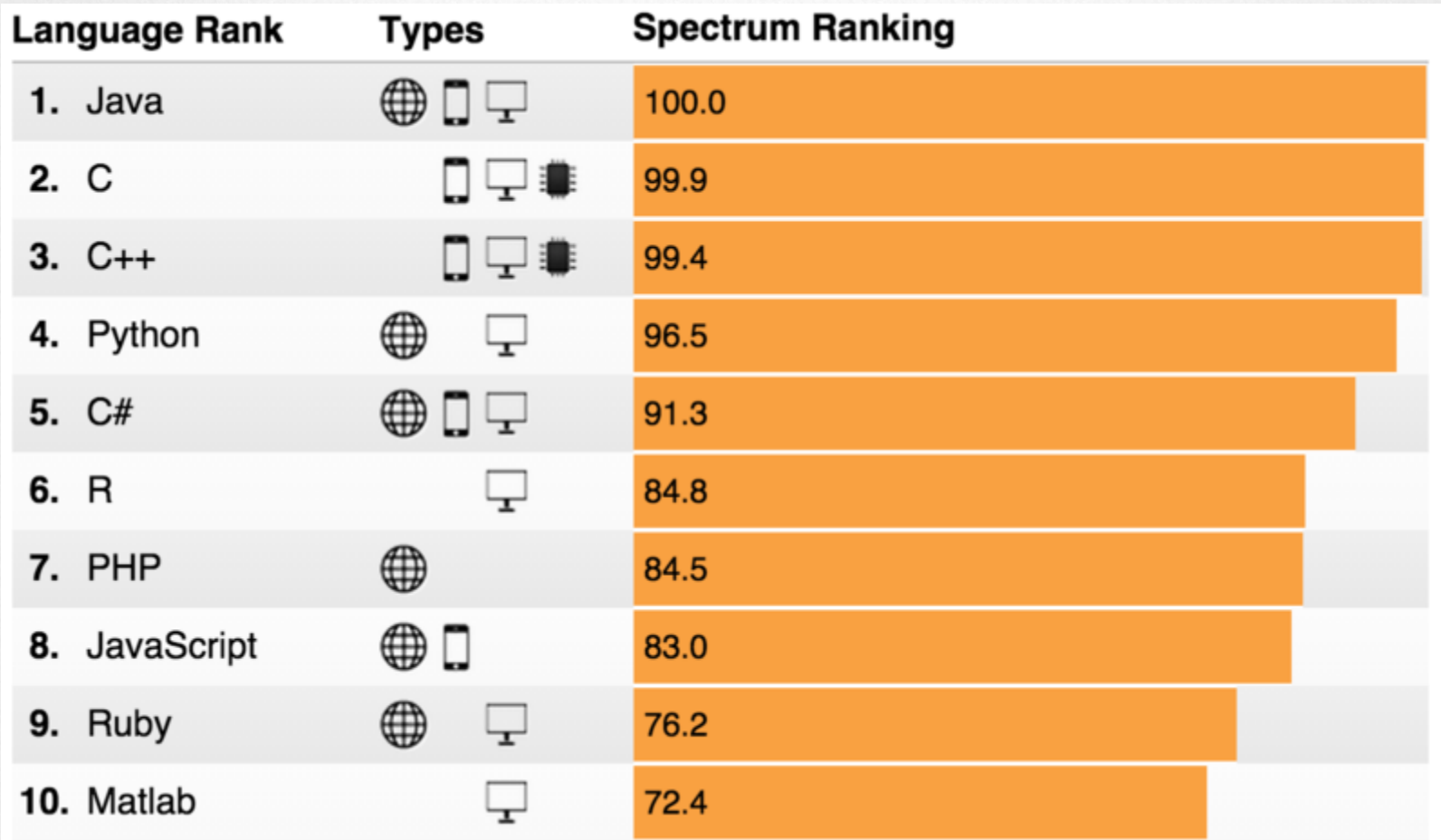


Popularity

2015-06-04



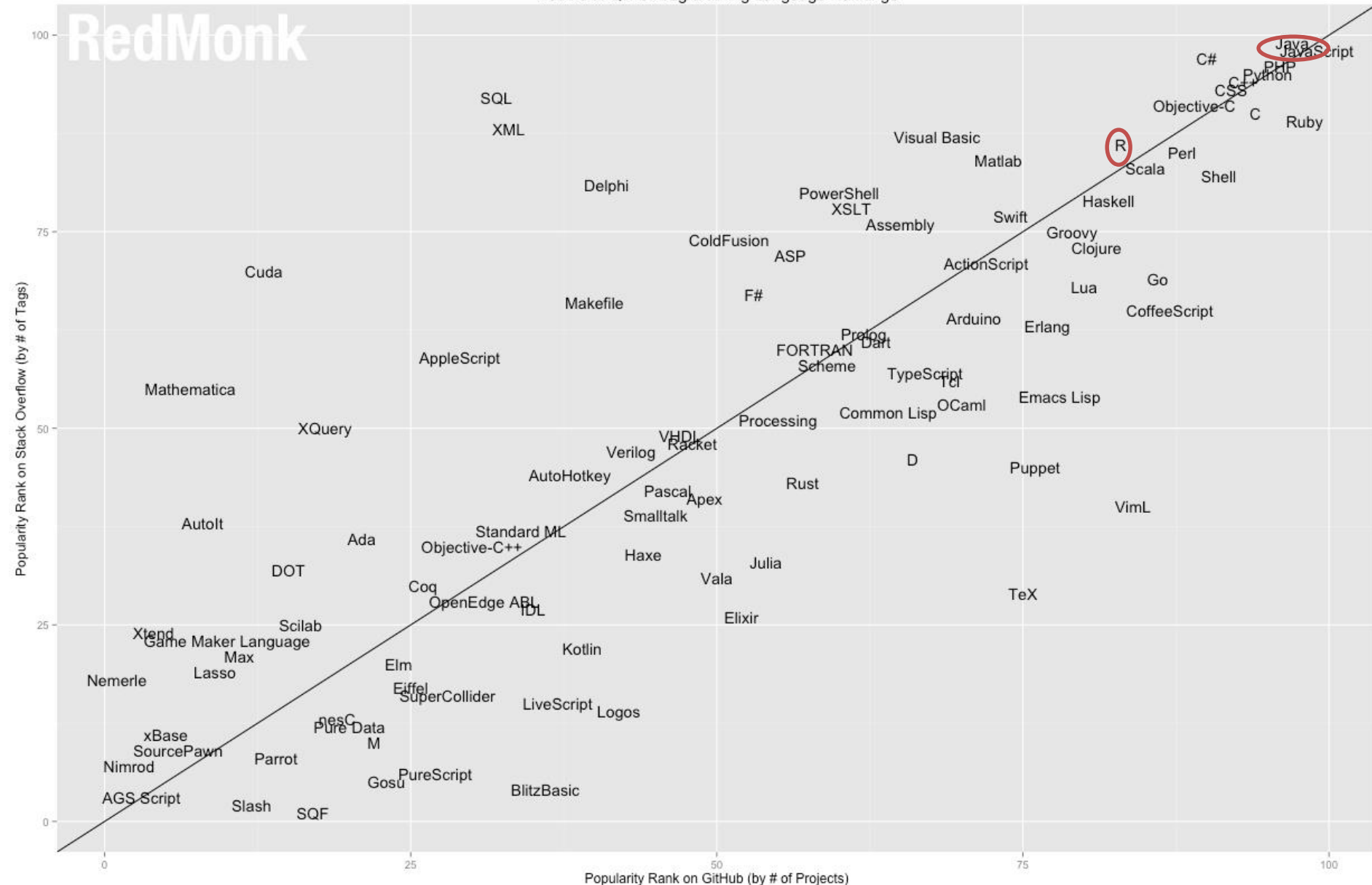
	PyPI	CRAN
Packages	60,806 Packages 35+ updates / day	6,727 package 20+ updates/day
Popularity (Tiobe)	6 th Rank, +0.67% 	12 th Rank, +1.06% 
 stackoverflow	430,604	93,943
 github SOCIAL CODING	549,014	87,306



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

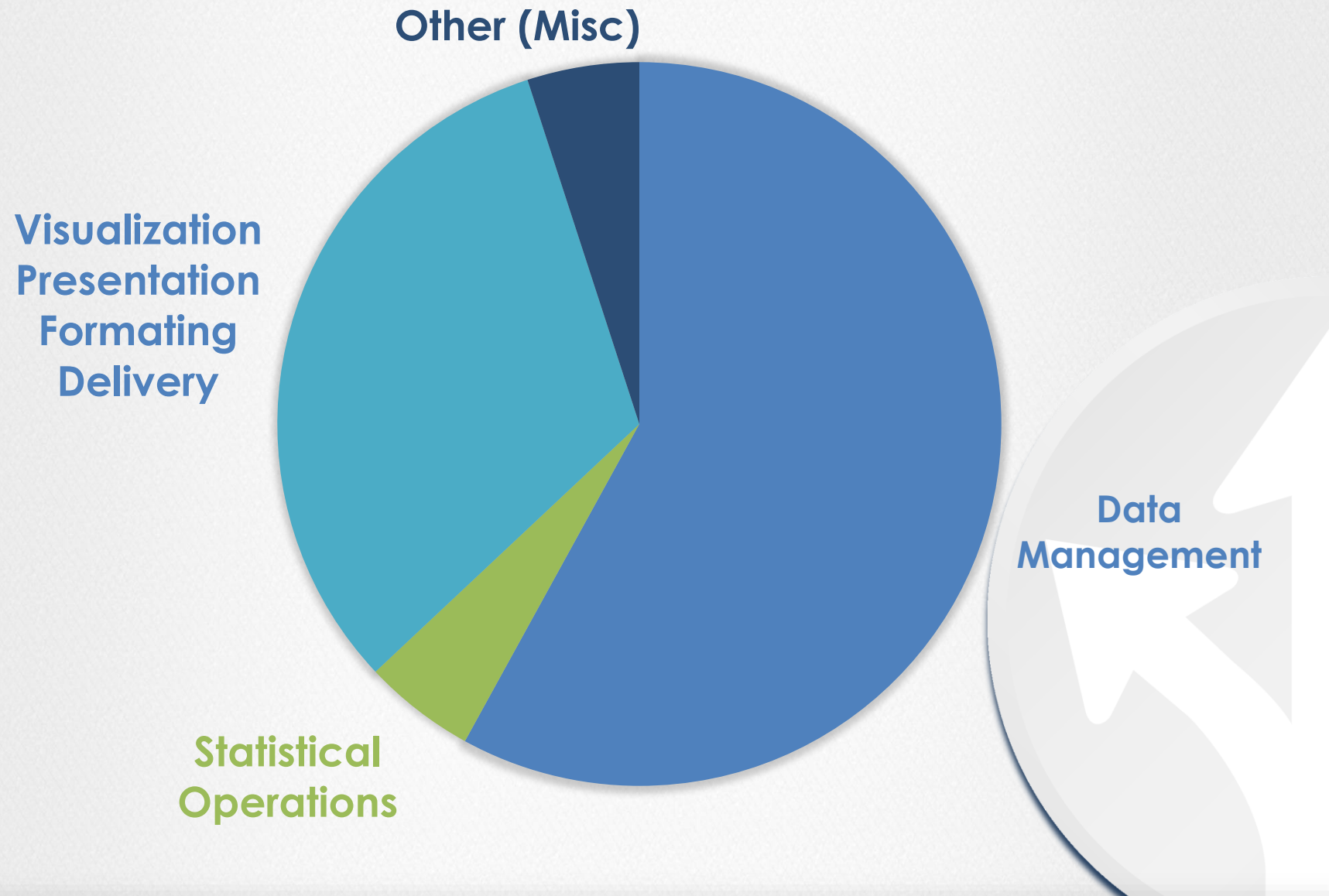
Domain Feature completeness

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 ...



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.

R Resources

ONLINE

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

Offline

The Art of R Programming

Norm Matloff

ISBN-13: 858-2592222227

ISBN-10: 1593273843

R in Action: Data Analysis and Graphics with R

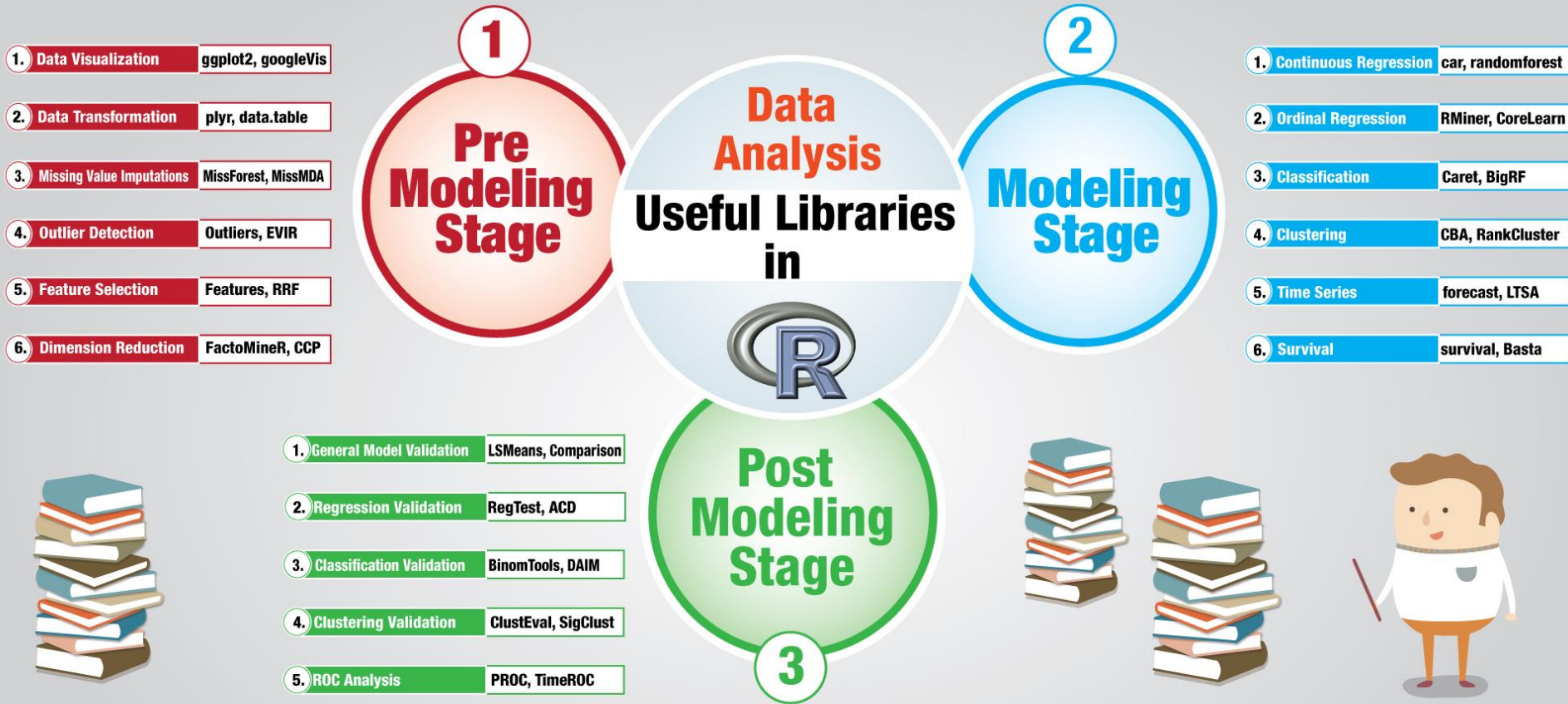
Robert Kabacoff

ISBN-13: 978-1617291388

ISBN-10: 1617291382



```
>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

R'S NOBELS



For Machine Learning

CARET

(CLASSIFICATION AND REGRESSION TRAINING)



Code readability, interactive programming

MAGRITTR
PIPER
BACKPIPE



Munging and data management

DATA.TABLES

DPLYR

<https://github.com/Rdatatable/data.table/wiki/Benchmarks-%3A-Grouping>



Simple Scaling-out

FOREACH ITERTOOLS



Increasing Performance

RCPP **GMATRIX**



Visualization

GGVIS **GGPLOT2**



APPENDIX



Given a vector of numbers (x), write a function (f) that returns a vector of numbers containing the *product* of every other number excluding the current index.

Example:

```
> x <- c( 1, 5, 2, 8 )
```

```
> f(x)
```

```
[1] 80 16 40 10
```

```
# 5*2*8, 1*2*8, 1*5*8, 1*2*5
```

Given a vector of numbers (x)
write a function (f) that returns a
vector of numbers containing the
product of every other integer
excluding the current index.

Example:

```
> x <- c( 1, 5, 2, 8 )  
> f(x)  
[1] 80 16 40 10  
  
# 5*2*8, 1*2*8, 1*5*8, 1*2*5
```

Solution:

```
f <- function(x) prod(x) / x
```


Write a function `f(x)` to accept an integer vector, and returns a vector with those numbers ... except for:

multiples of 3 = "Fizz"

multiples of 5 = "Buzz".

multiples of 3 and 5 = "FizzBuzz"

Example:

```
> x <- 1:20
```

```
> f(1:20)
```

```
1 2 Fizz 4 Buzz Fizz 7 8 Fizz Buzz 11  
Fizz 13 14 FizzBuzz 16 17 Fizz 19 Buzz
```

Solution:

```
f <- function(x)  
  ifelse( x %% 15 == 0, "FizzBuzz",  
    ifelse( x %% 3 == 0, "Fizz",  
      ifelse( x %% 5 == 0, "Buzz", x ) ) )
```

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?