

Tools for Data Science

Vadim Y. Bichutskiy
@vybstat

Data Science Seminar, GMU
April 10, 2015

So you want to be a data scientist?

- Good news
 - Data is everywhere
 - “Big Data”, “Analytics”, “Data Science” is changing the world
 - Hot and sexy
 - Lots of opportunity to get creative and innovate
 - Many open problems
 - Fun!
 - Demand is off the charts / low supply
 - High salaries
- Bad news
 - Requires lots of education: PhD is NOT enough
 - Can be overwhelming and stressful
 - Theory, practical tools, experience
 - Long working hours
 - Not enough sleep
 - Bad for health?
 - Versatile, flexible, curious
 - Continuous training



BRIEFING ROOM

ISSUES

THE ADMINISTRATION

PARTICIPATE

1600 PENN

Search



Home • The White House Blog

The White House Blog

[Subscribe](#)

Our Top Stories



[Another Step Toward Equality for LGBT Workers](#)



[Weekly Address: Reaching a Comprehensive and Long-Term Deal on Iran's Nuclear Program](#)



[Continuing Our Focus on Solar Energy](#)



[Innovative Job-Training Programs Are Important. Here's Why:](#)

The White House Names Dr. DJ Patil as the First U.S. Chief Data Scientist



Megan Smith

February 18, 2015
04:48 PM EDT

Share This Post



Today, I am excited to welcome Dr. DJ Patil as Deputy Chief Technology Officer for Data Policy and Chief Data Scientist here at the White House in the Office of Science and Technology Policy. President Obama has prioritized bringing top technical talent like DJ into the federal government to harness the power of technology and innovation to help government better serve the American people.

Across our great nation, we've begun to see an acceleration of the power of data to deliver value. From early open data work by the National Oceanic and Atmospheric Administration (NOAA), which provides data that enables weather forecasts to come directly to our mobile phones, to powering GPS systems that feed geospatial data to countless apps and services — government data has supported a transformation in the way we live today for the better.

DJ joins the White House following an incredible career as a data scientist — a term he helped coin — in the public and private sectors, and in academia. Most recently, DJ served as the Vice

[Technology](#)

WHITEHOUSE.GOV IN YOUR INBOX

Sign up for email updates from President Obama and Senior Administration Officials

Your Email Address

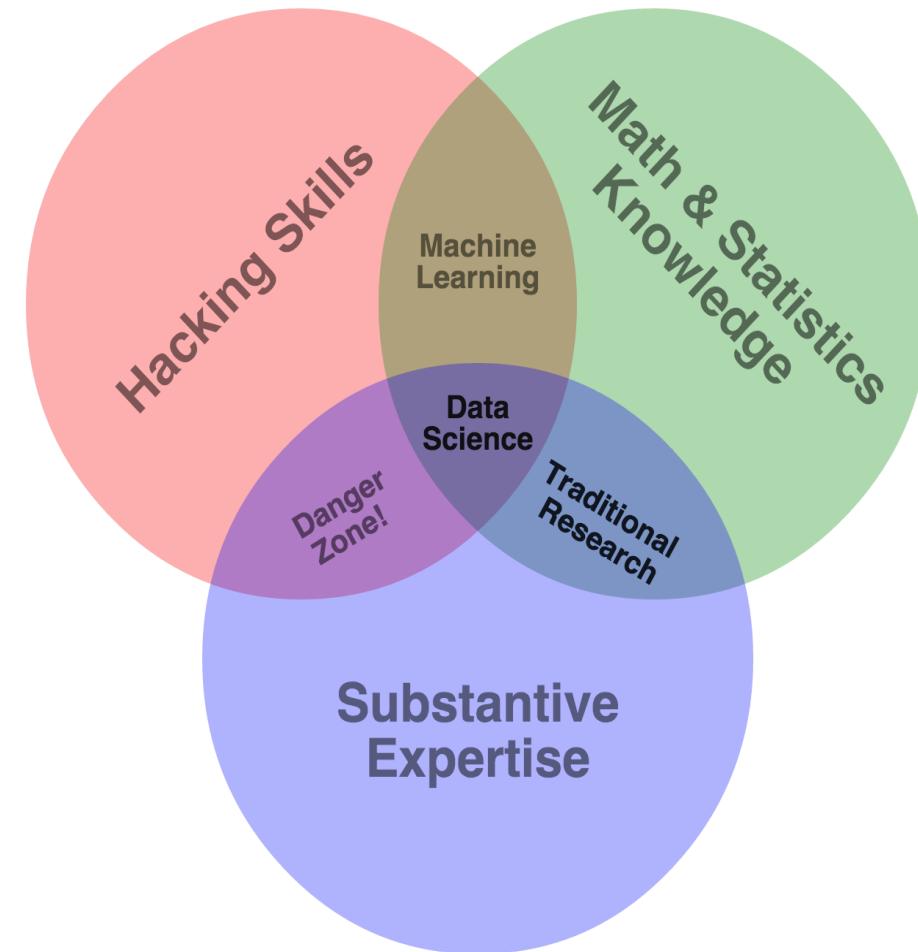
Submit

PHOTOS OF THE DAY

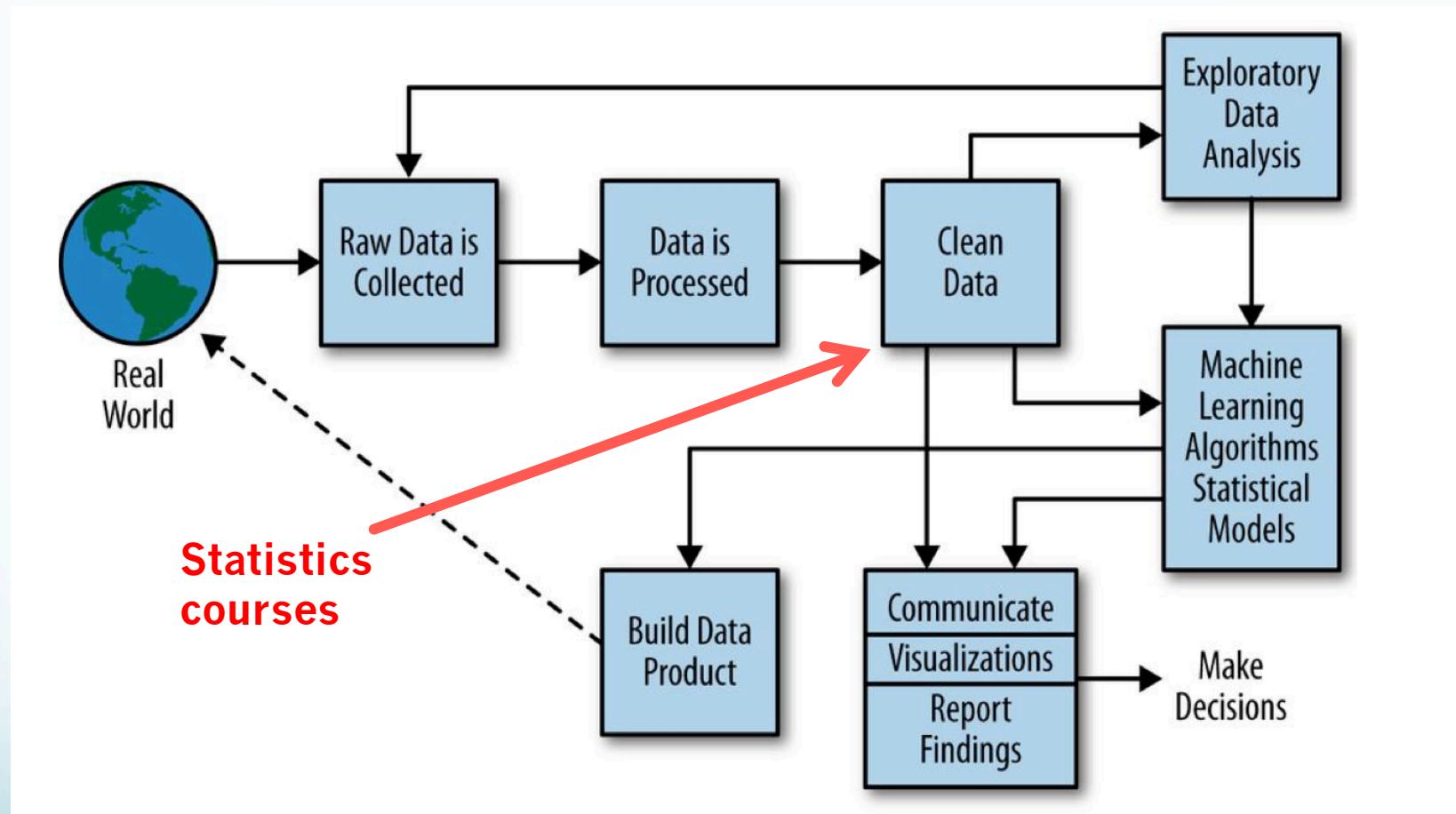


What's Data Science?

The Data Science Venn Diagram



Data scientists: “Create order from chaos”



O’Neil, Cathy and Schutt, Rachel, *Doing Data Science: Straight Talk from the Frontline*, O’Reilly, 2014

Data collection, processing, cleaning is 80% of the effort

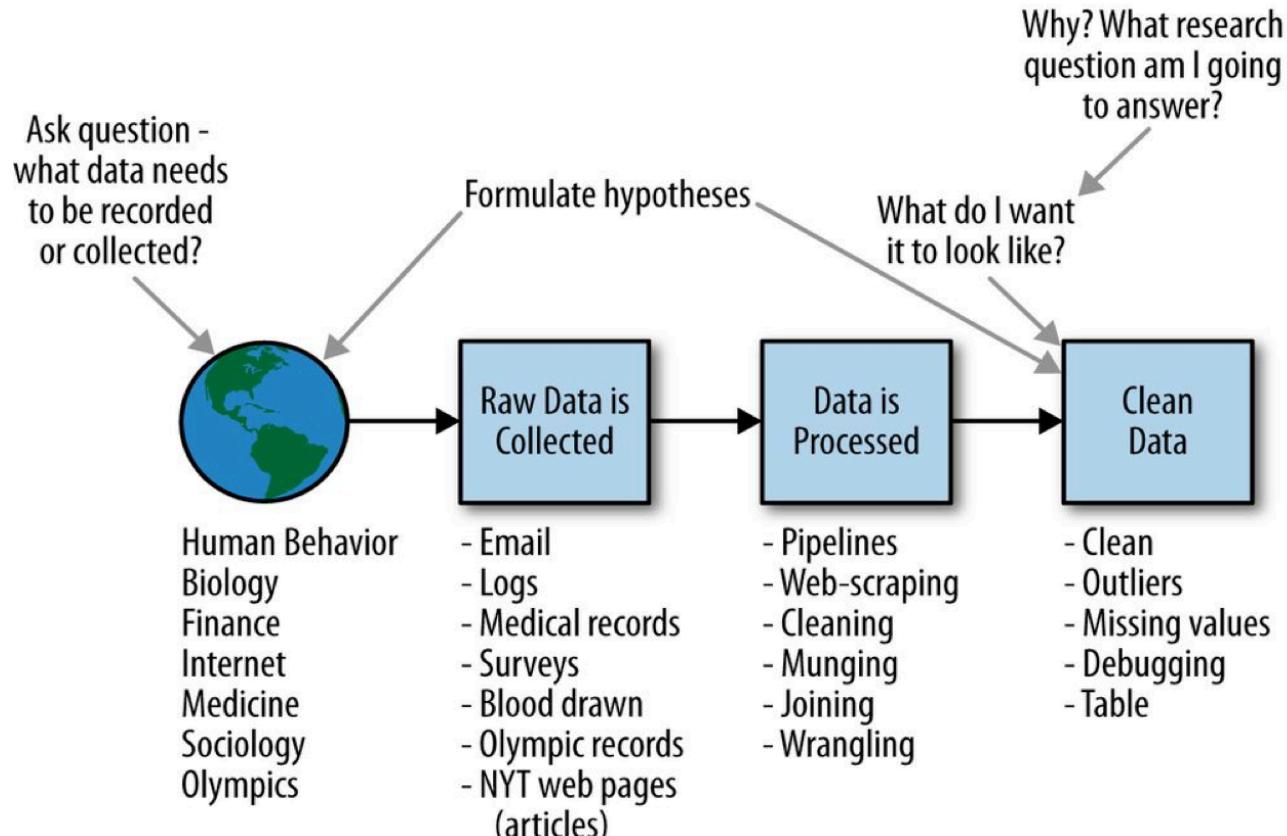
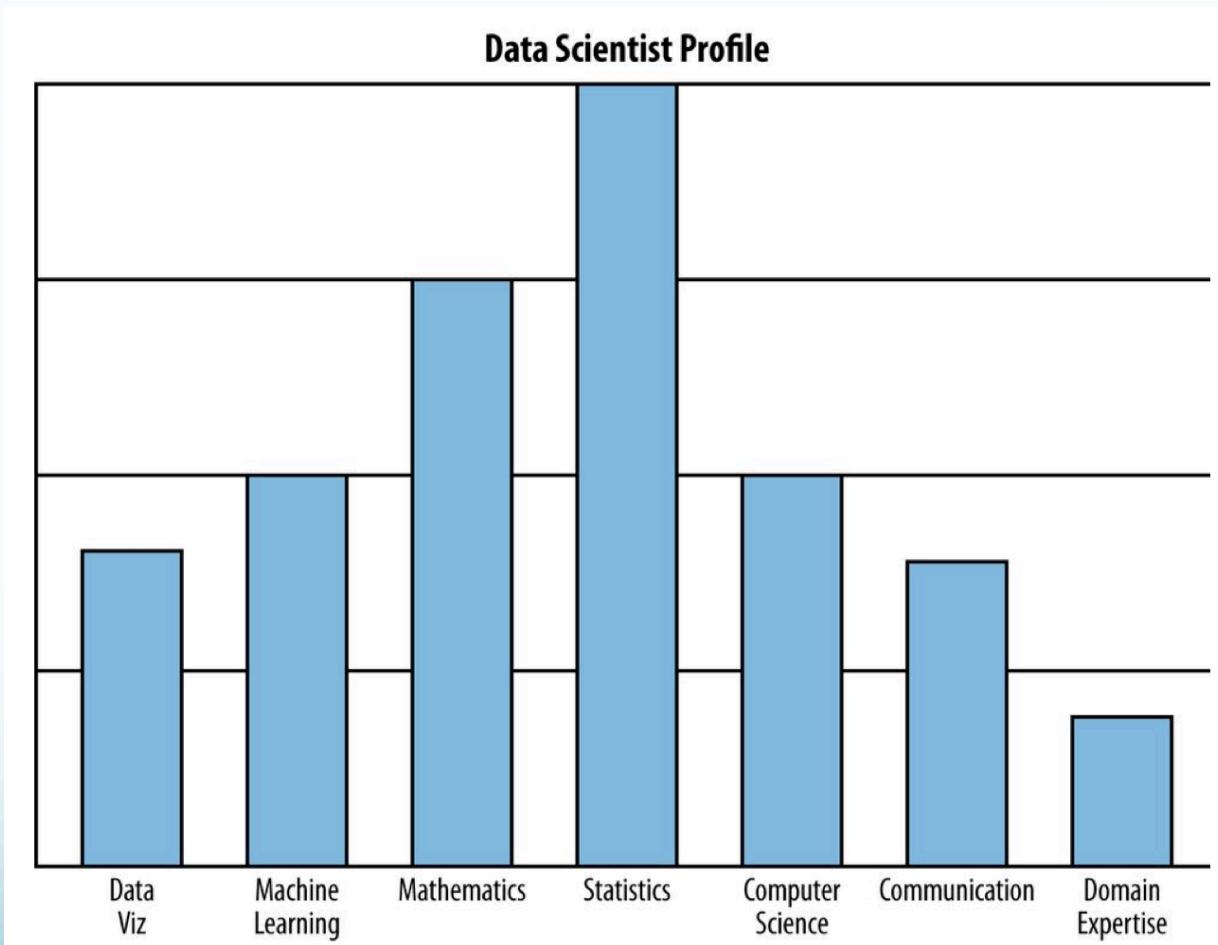


Figure 2-3. The data scientist is involved in every part of this process

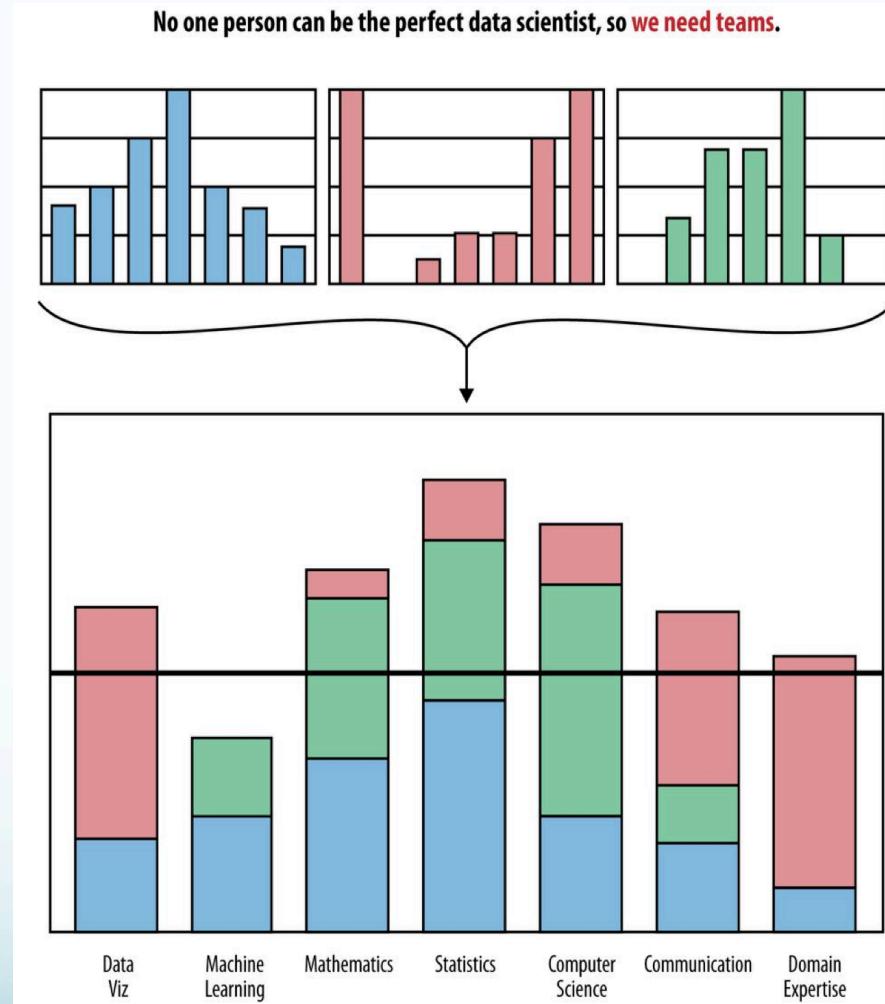
O'Neil, Cathy and Schutt, Rachel, *Doing Data Science: Straight Talk from the Frontline*, O'Reilly, 2014

Stats/CSI PhD

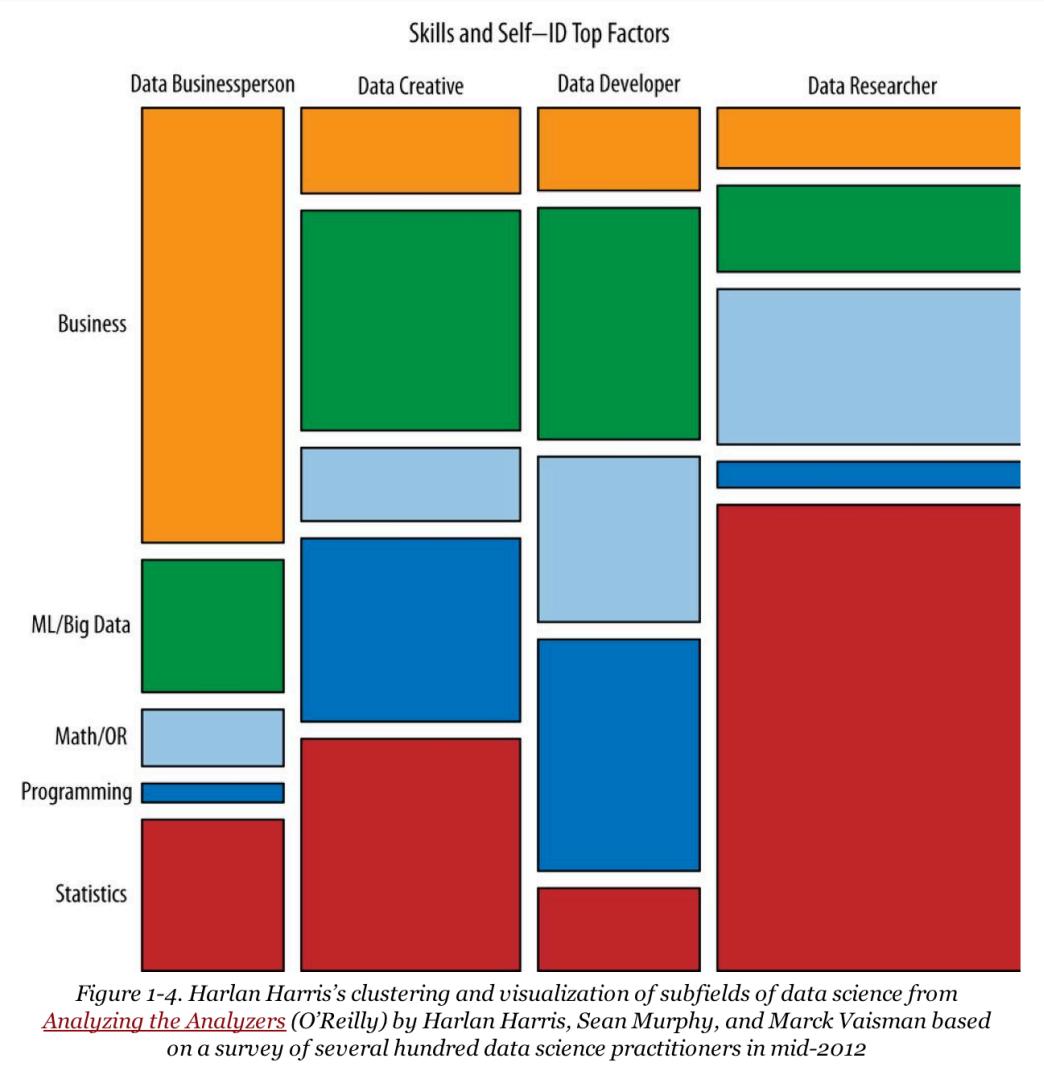


O'Neil, Cathy and Schutt, Rachel, *Doing Data Science: Straight Talk from the Frontline*, O'Reilly, 2014

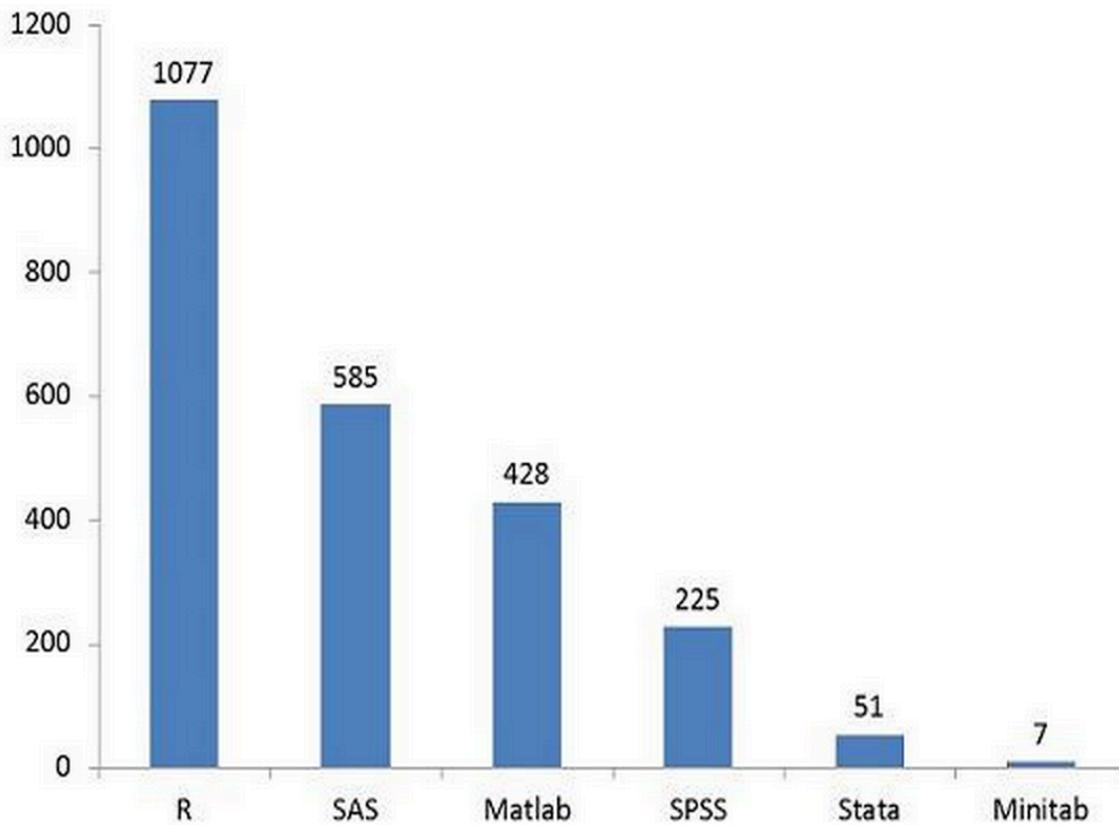
“Data science is a team sport” --DJ Patil



O'Neil, Cathy and Schutt, Rachel, *Doing Data Science: Straight Talk from the Frontline*, O'Reilly, 2014



Statistical Analysis Tools Listed in Data Scientist Job Descriptions



Kirk Borne @KirkDBorne · Apr 3

Tech Skills that #DataScientist Jobs Require: bit.ly/1DHBGbe #abdsc
#BigData #Python #Rstats @DataScienceCtrl



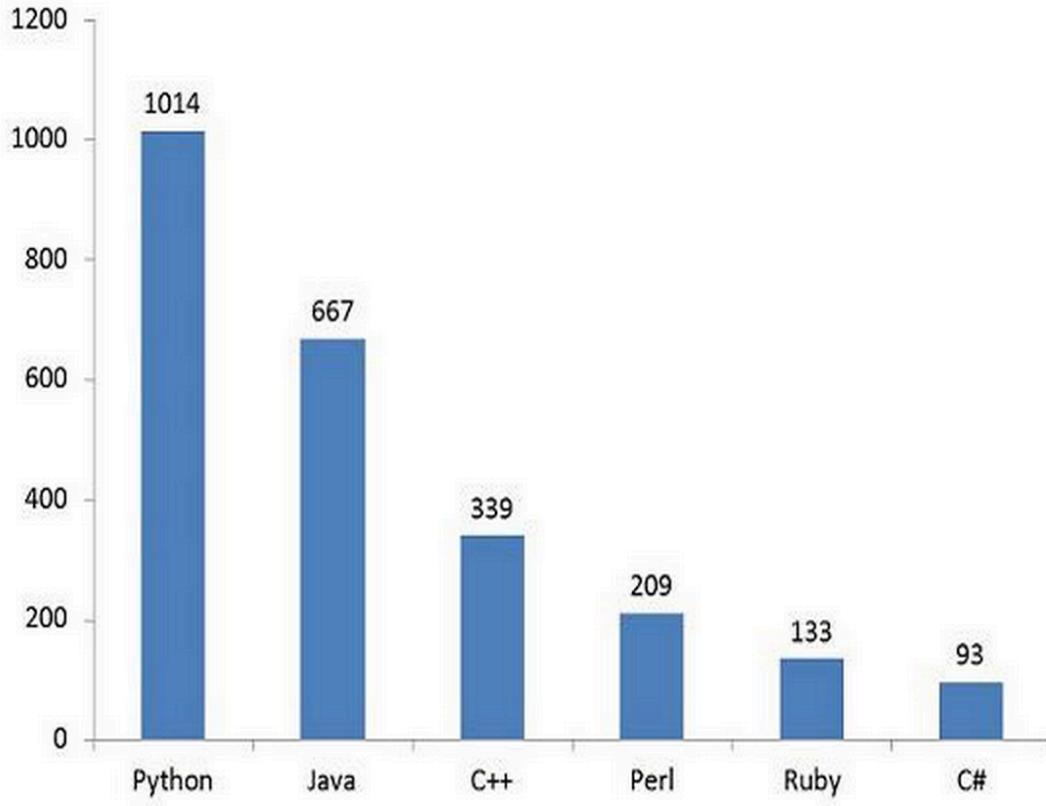
72



54

...

Programming Languages Listed in Data Scientist Job Descriptions



Kirk Borne @KirkDBorne · Apr 3

Tech Skills that #DataScientist Jobs Require: bit.ly/1DHBGbe #abdsc
#BigData #Python #Rstats @DataScienceCtrl



72



54

...

Data Science Skills

- Core
 - R
 - Python
 - SQL/NoSQL/database concepts
 - Unix command line
 - Statistics/machine learning/CS
 - Graph Theory/Networks
 - Data visualization/dashboards: Tableau, D3
 - Data representation: JSON, XML
 - Communication, domain expertise
- Project/position dependent
 - Java/C++
 - Amazon Web Services/Cloud computing
 - Hadoop
 - JavaScript/PHP/Web frameworks
- Emerging
 - Scala
 - Swift
 - Spark/Cluster computing
 - Real-time Analytics
 - Docker, Vagrant

Tools Usage

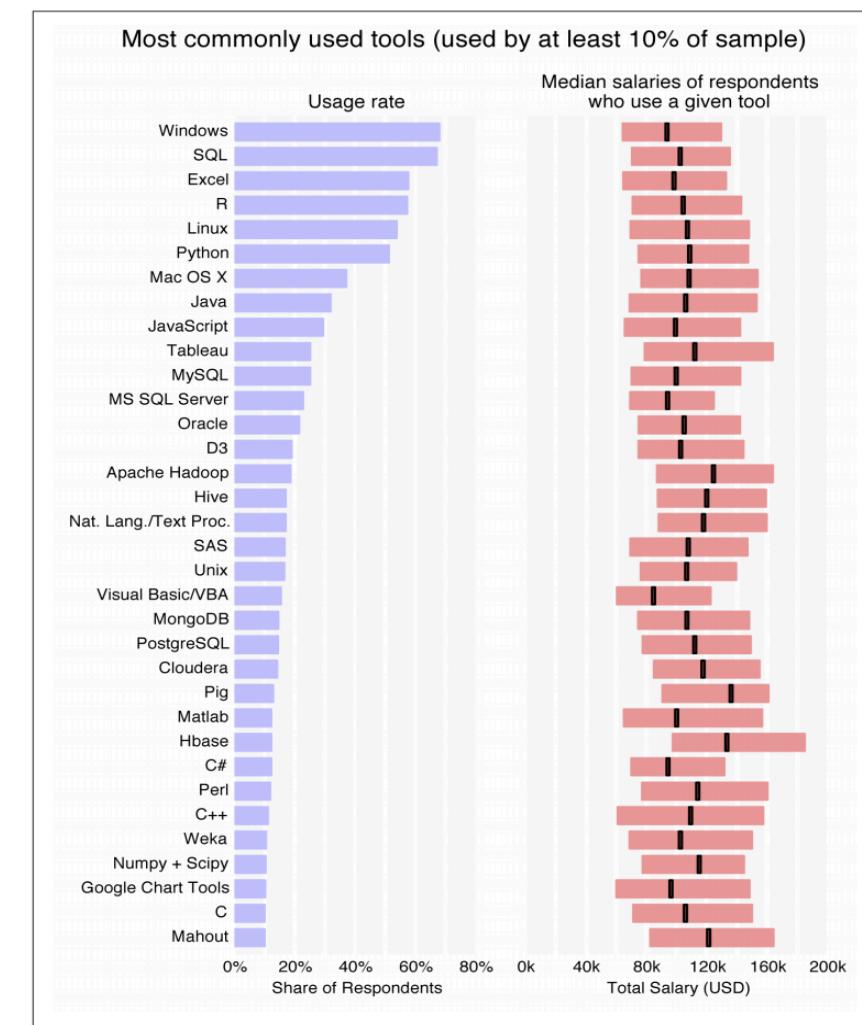


Figure 1-10. Most commonly used tools

Tool Salaries

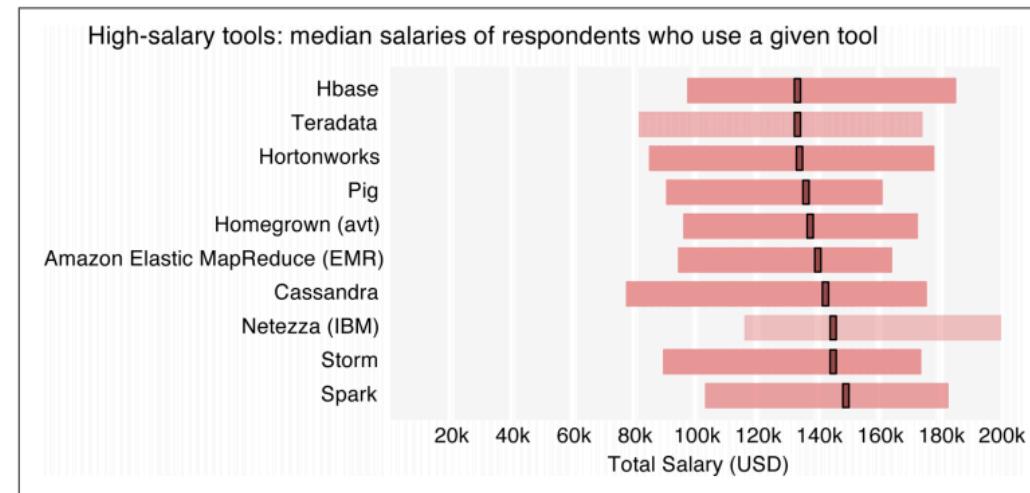


Figure 1-11. High-salary tools: median salaries of respondents who use a given tool

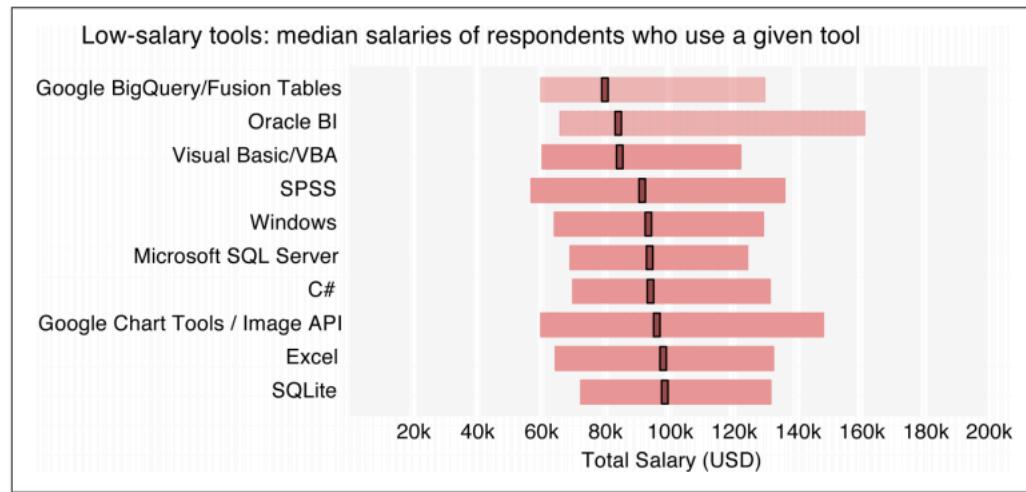
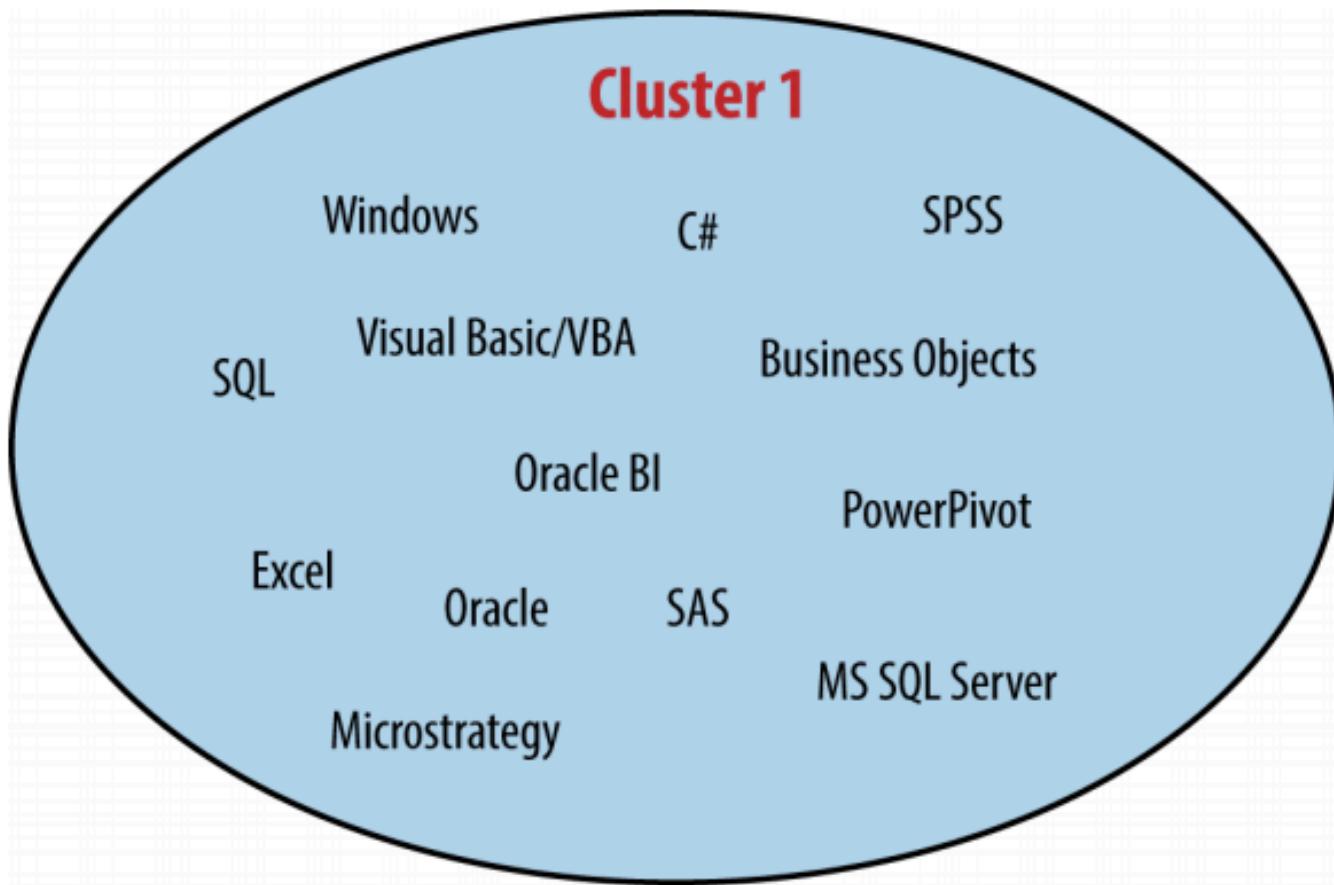
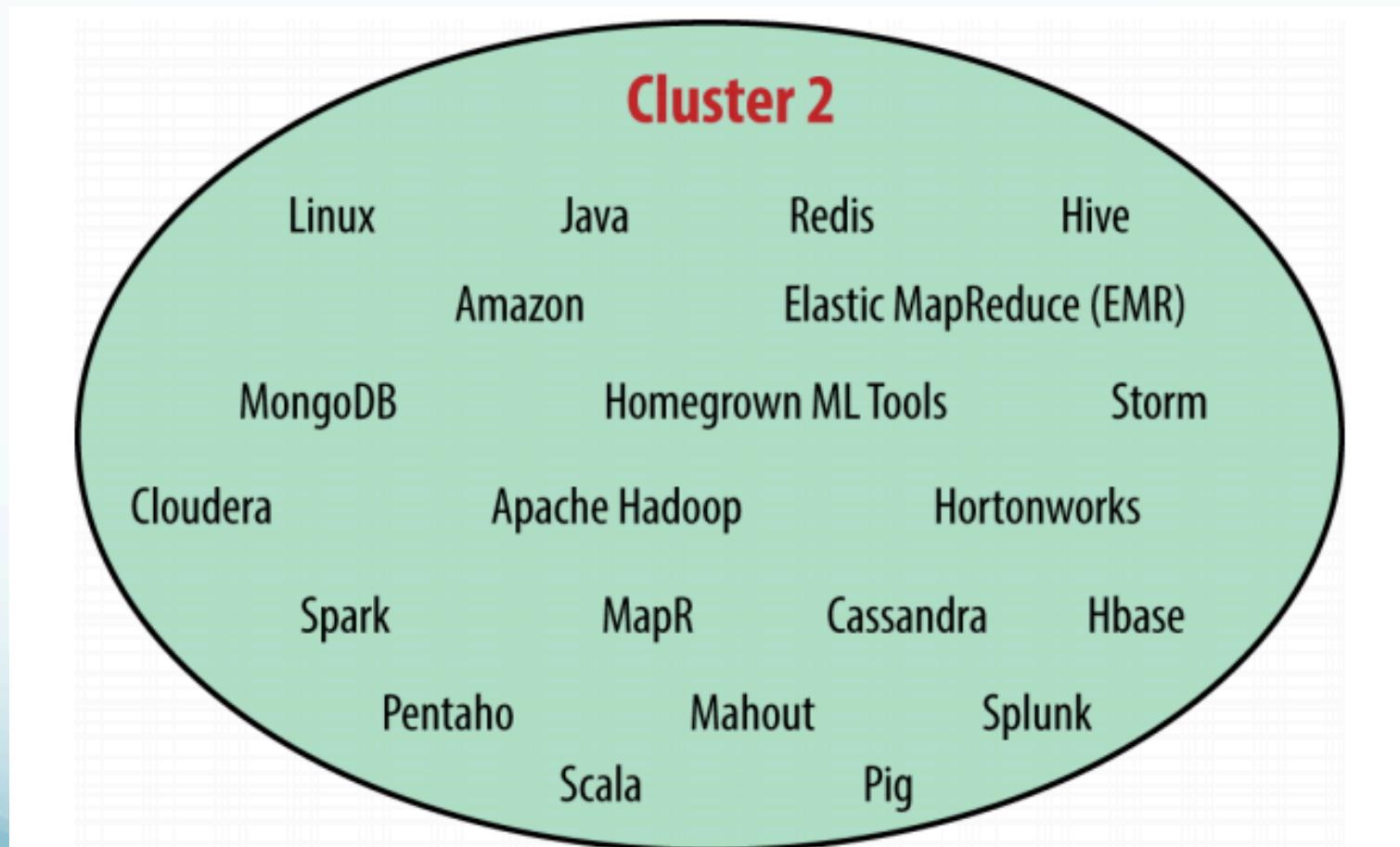


Figure 1-12. Low-salary tools: median salaries of respondents who use a given tool

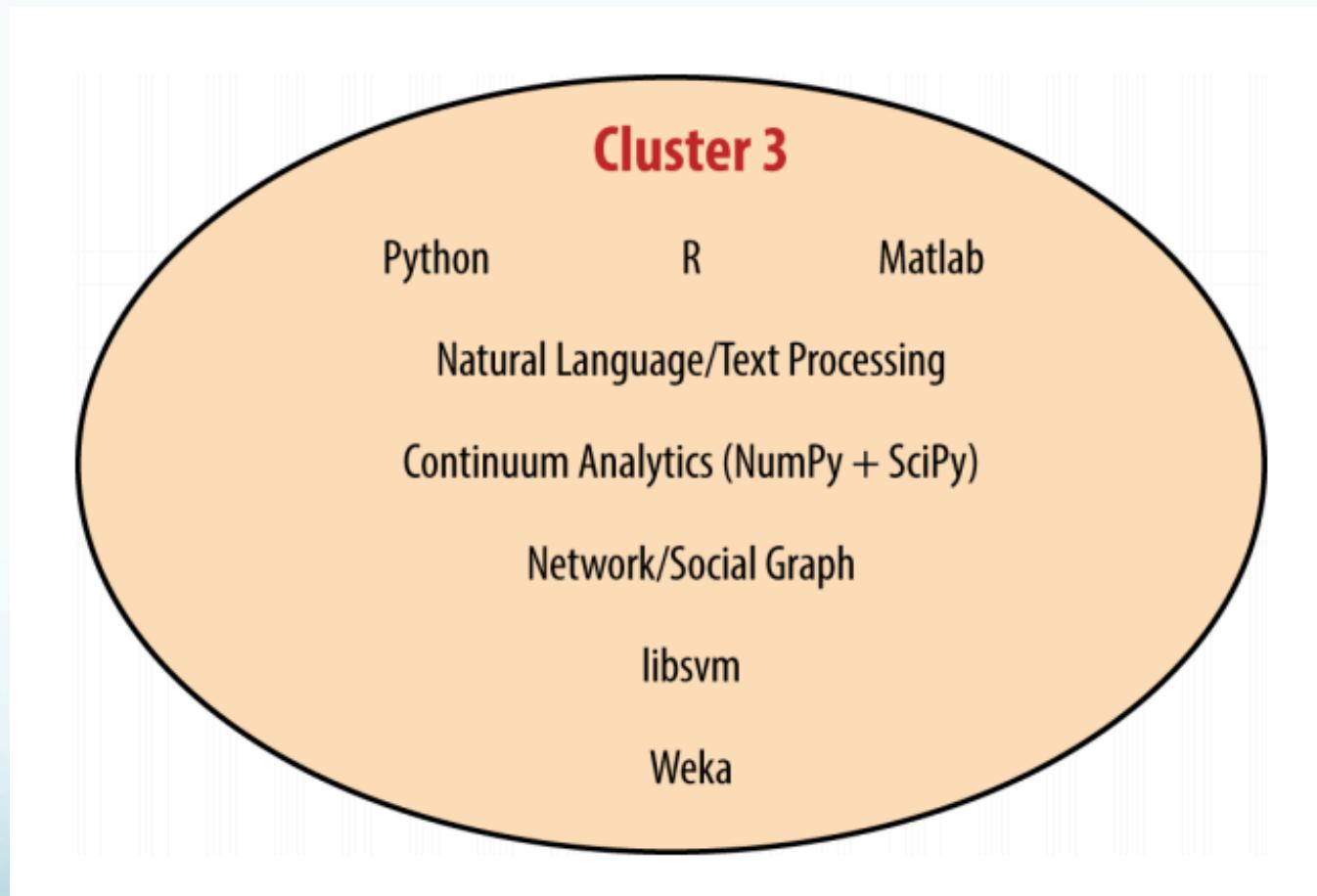
“Microsoft-Excel-SQL”



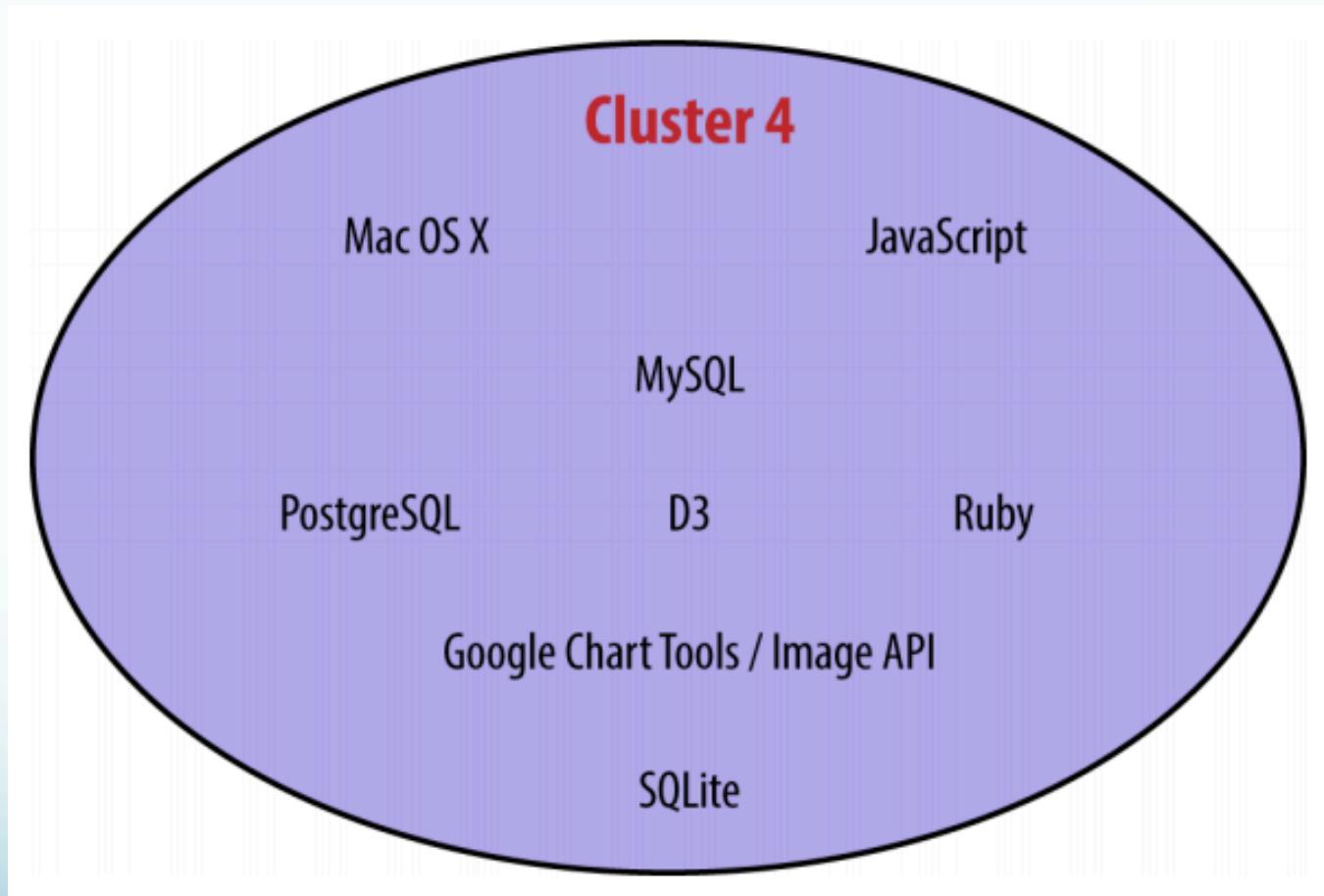
“Hadoop-Java-Cloud Computing”



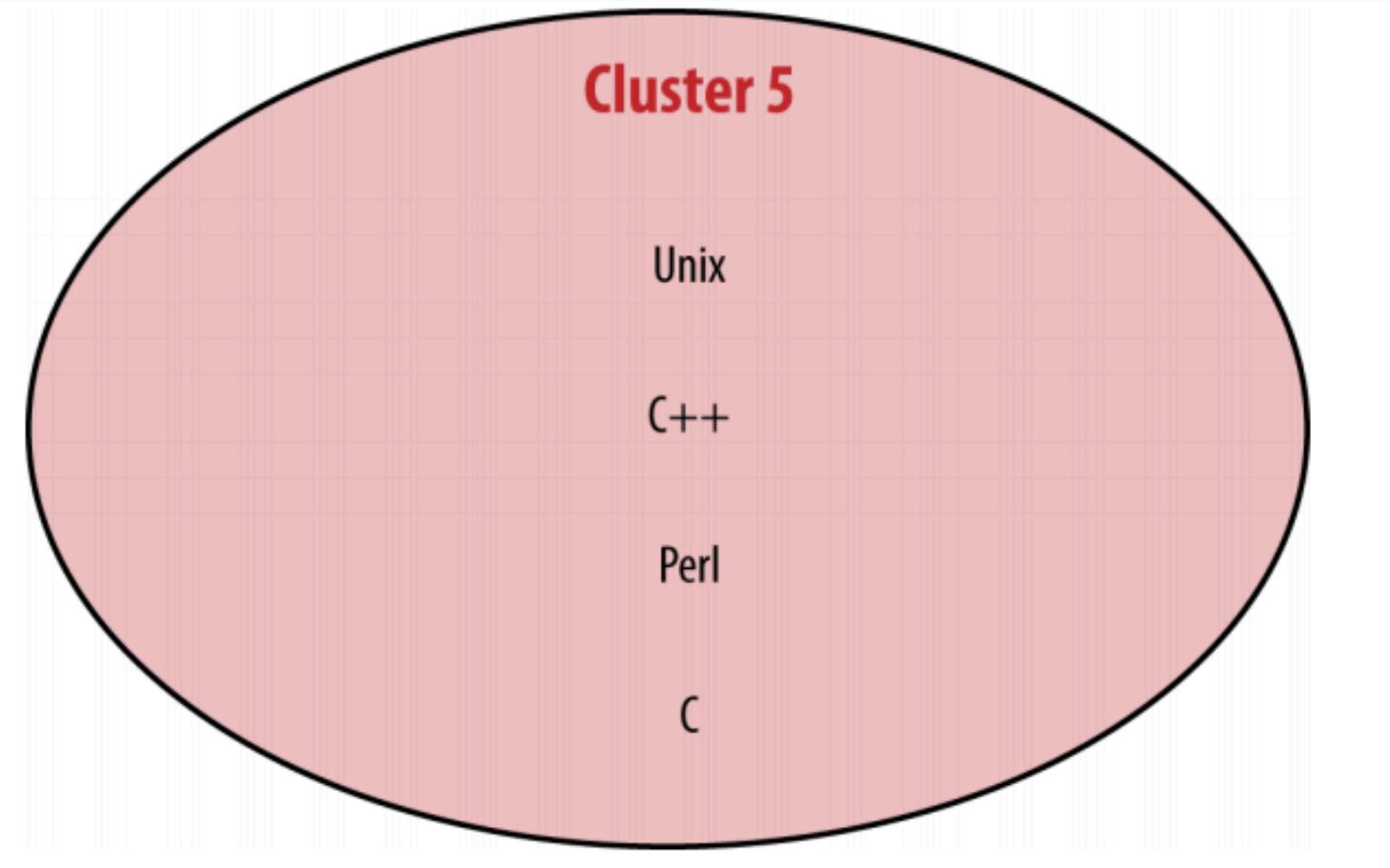
“R-Python-Analytics”



“MySQL-D3-JavaScript”



“Old tools”



Amazon MLaaS

AWS Official Blog

Amazon Machine Learning – Make Data-Driven Decisions at Scale

by Jeff Barr | on 09 APR 2015 | in [Amazon Machine Learning](#) | [Permalink](#)

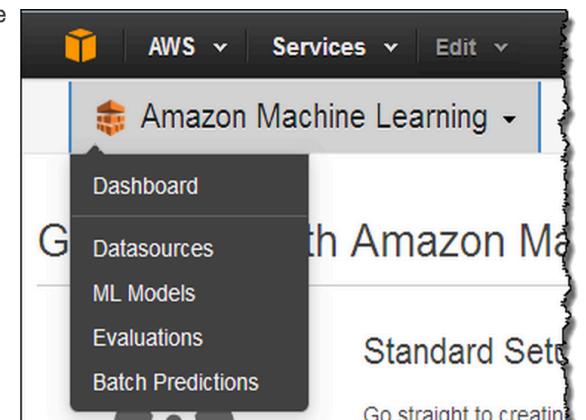
Today, it is relatively straightforward and inexpensive to observe and collect vast amounts of operational data about a system, product, or process. Not surprisingly, there can be tremendous amounts of information buried within gigabytes of customer purchase data, web site navigation trails, or responses to email campaigns. The good news is that all of this data can, when properly analyzed, lead to statistically significant results that can be used to make high-quality decisions. The bad news is that you need to find data scientists with relevant expertise in machine learning, hope that your infrastructure is able to support their chosen tool set, and hope (again) that the tool set is sufficiently reliable and scalable for production use.

The science of Machine Learning (often abbreviated as ML) provides the mathematical underpinnings needed to run the analysis and to make sense of the results. It can help you to turn all of that data into high-quality predictions by finding and codifying patterns and relationships within it. Properly used, Machine Learning can serve as the basis for systems that perform fraud detection (is this transaction legitimate or not?), demand forecasting (how many widgets can we expect to sell?), ad targeting (which ads should be shown to which users?), and so forth.

Introducing Amazon Machine Learning

Today we are introducing [Amazon Machine Learning](#). This new AWS service helps you to use all of that data you've been collecting to improve the quality of your decisions. You can build and fine-tune predictive models using large amounts of data, and then use [Amazon Machine Learning](#) to make predictions (in batch mode or in real-time) at scale. You can benefit from machine learning even if you don't have an advanced degree in statistics or the desire to setup, run, and maintain your own processing and storage infrastructure.

I'll get to the details in just a minute. Before I do so, I'd like to review some of the terminology and the concepts that you need to know in order to fully understand what machine learning does and how you can take advantage





Vadim Y. Bichutskiy

@vybstat

.@Amazon launches machine learning service. Great, but doesn't replace data scientists Statisticians maybe? #BigData
aws.amazon.com/blogs/aws/amaz...



RETWEETS

2

FAVORITES

4



7:59 PM - 9 Apr 2015



Vadim Y. Bichutskiy

@vybstat

Training models is the easiest part of
#BigData #Analytics #DataScience
aws.amazon.com/blogs/aws/amaz...



RETWEET

1

FAVORITES

3



8:02 PM - 9 Apr 2015

Resources (1)

- R
 - <http://www.r-project.org/>
 - <http://www.rstudio.com/>
- Python
 - <https://www.python.org/>
- JSON
 - <http://json.org/>
- Amazon Web Services
 - <http://aws.amazon.com/>
 - <http://aws.amazon.com/blogs/aws/>
- Hadoop
 - <https://hadoop.apache.org/>

Resources (2)

- Scala
 - <http://www.scala-lang.org/>
- Spark
 - <https://spark.apache.org/>
 - <https://spark.apache.org/docs/latest/>
- Docker
 - <https://www.docker.com/>
- Vagrant
 - <https://www.vagrantup.com/>
- Swift
 - apple.co/1CAAQQA