

From Zero to Hero

A Hands-On Approach to R Stats

Pierre Lafortune
pierre.a.lafortune@gmail.com

Agenda

- R Introduction
- Common Use cases
- Problem Statement
- Live Coding
- Participant Challenges

R Quick Launch

- Statistical software environment
- Built on S programming language
- Named R after originators Robert Gentleman and Ross Ihaka
- Created in 1993
- Open-source under General Public License
- Source-code written in C, Fortran, and R
- Over 6,000 packages added by contributors



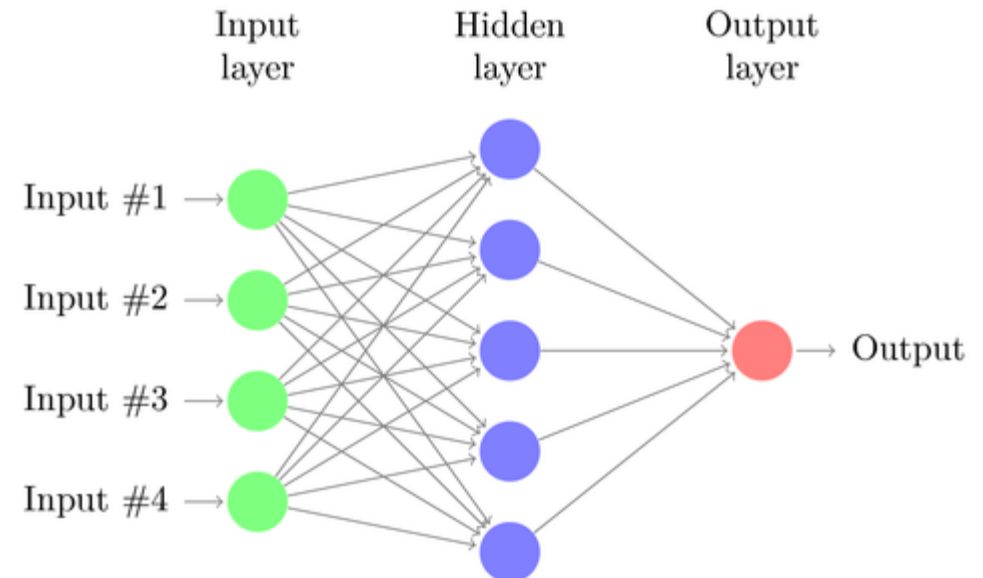
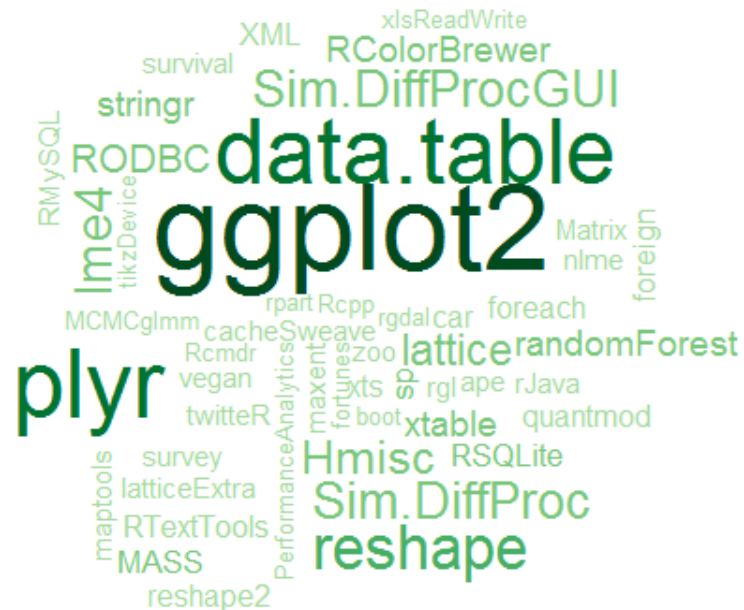
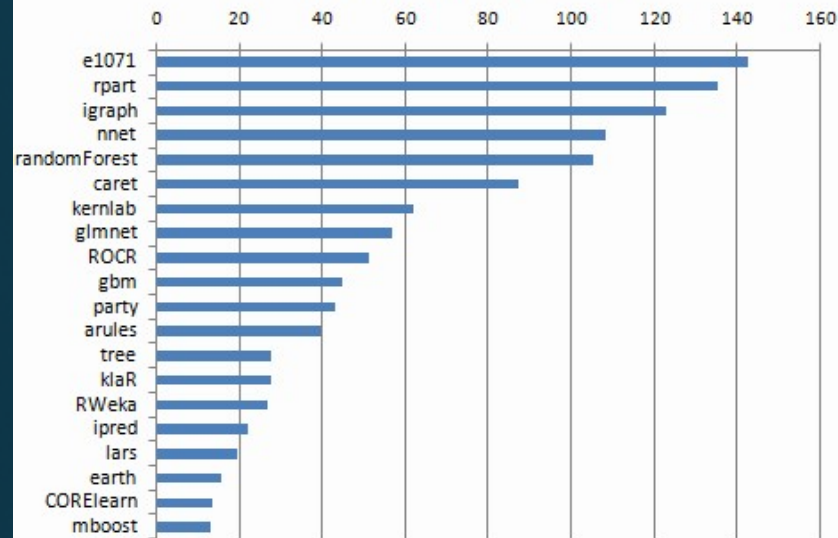
They use(d) R

<http://www.revolutionanalytics.com/companies-using-r>

<http://www.rstudio.com/>



Top 20 R Machine Learning packages, by Downloads (000) from CRAN

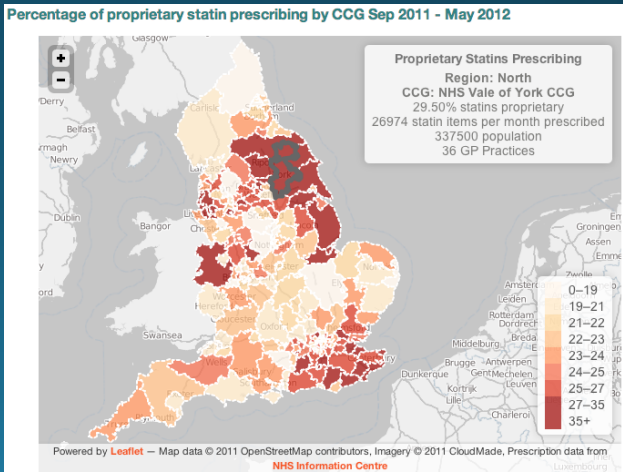
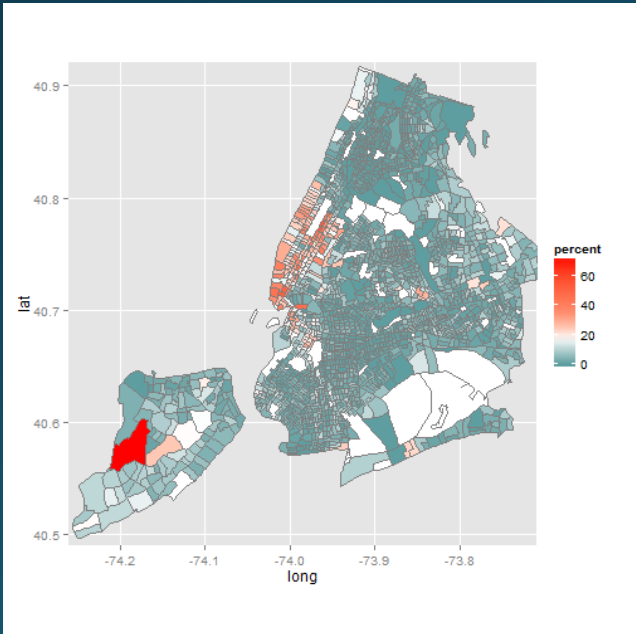
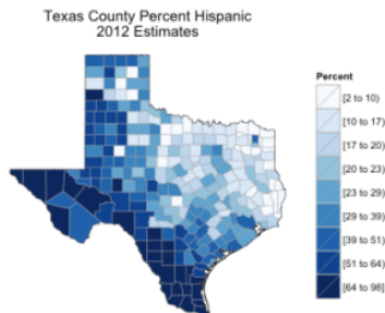




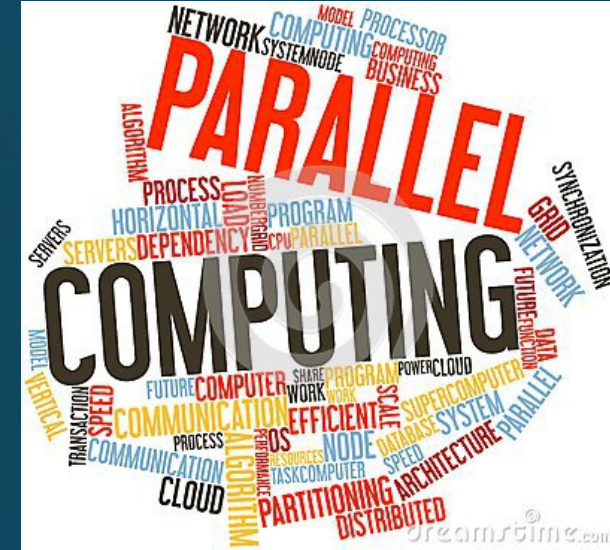
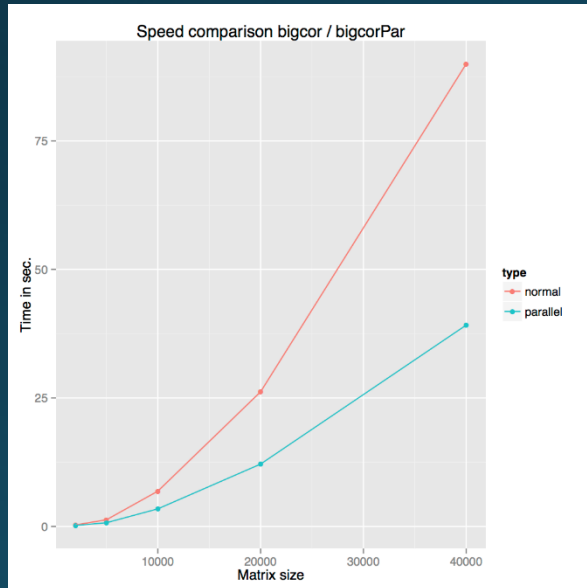
Learn to Map Census Data in R

JUNE 25, 2015 / ARI LAMSTEIN / 6 COMMENTS

Today I am happy to announce a new free email course: **Mapping Census Data in R**. You can sign up via the form at the bottom of this post. The course is designed to provide similar information to what I covered in my tutorial [Analyzing US Census Data with R](#). In short, it will teach you how to create choropleth maps of US demographics such as this¹:

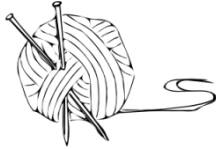


Parallel Processing





knitr
Elegant, flexible and fast
dynamic report generation with R

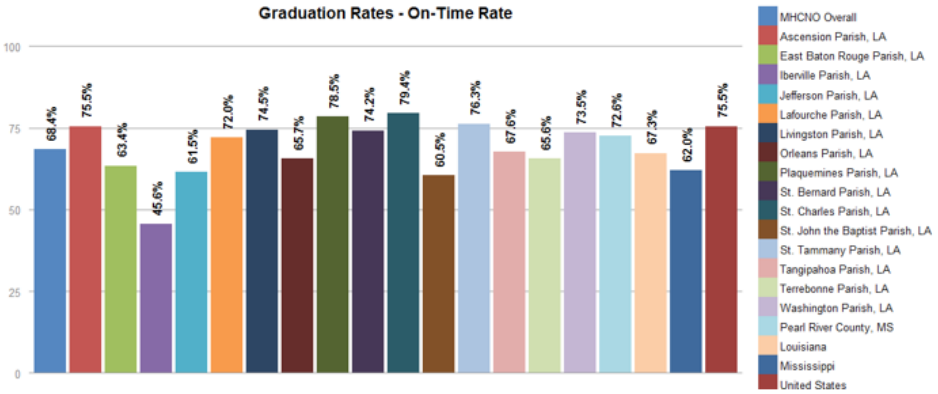
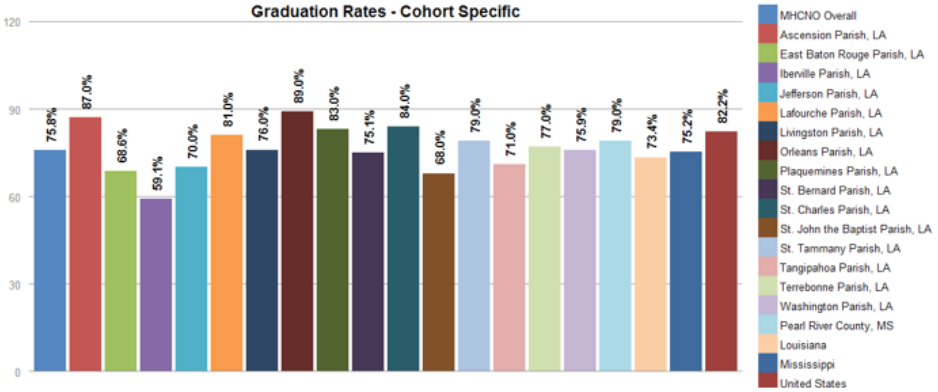


Hispanic Population by Race Alone, Percent

Report Area	White	Black	Asian	Native American / Alaska Native	Native Hawaiian / Pacific Islander	Some Other Race	Multiple Races
Report Area	75.9%	4.86%	0.43%	0.52%	0.1%	11.16%	7.04%
Calhoun County, FL	75.77%	0%	0%	0%	0%	17.31%	6.92%
Charlotte County, FL	86.85%	2.17%	0.3%	0.7%	0%	6.49%	3.49%
Duval County, FL	74.42%	5.32%	0.45%	0.5%	0.12%	11.63%	7.56%
Holmes County, FL	62.77%	2.14%	0%	0%	0%	30.02%	5.07%
Florida	83.87%	3%	0.19%	0.38%	0.03%	9.58%	2.96%
United States	65.37%	2.08%	0.33%	0.91%	0.08%	26.65%	4.58%

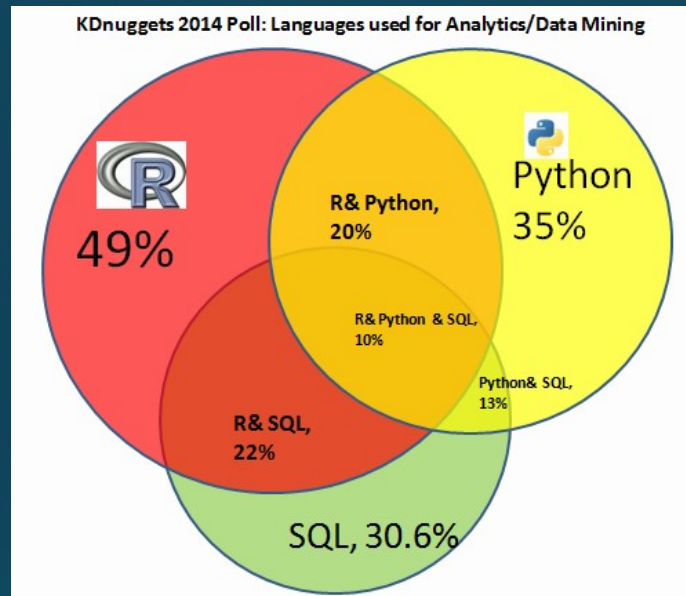
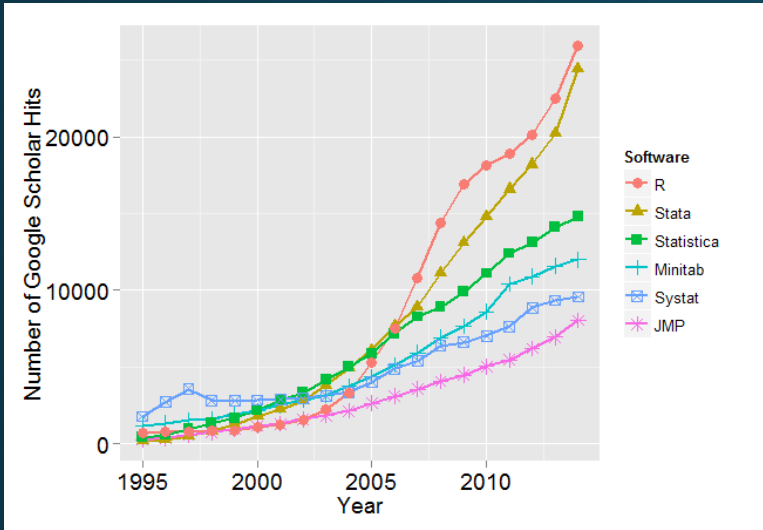
Non-Hispanic Population by Race Alone, Total

Report Area	White	Black	Asian	Native American / Alaska Native	Native Hawaiian / Pacific Islander	Some Other Race	Multiple Races
Report Area	657,540	268,072	39,124	2,631	812	2,014	24,957
Calhoun County, FL	11,334	1,947	92	219	0	0	256
Charlotte County, FL	138,848	9,201	2,000	308	5	217	2,285

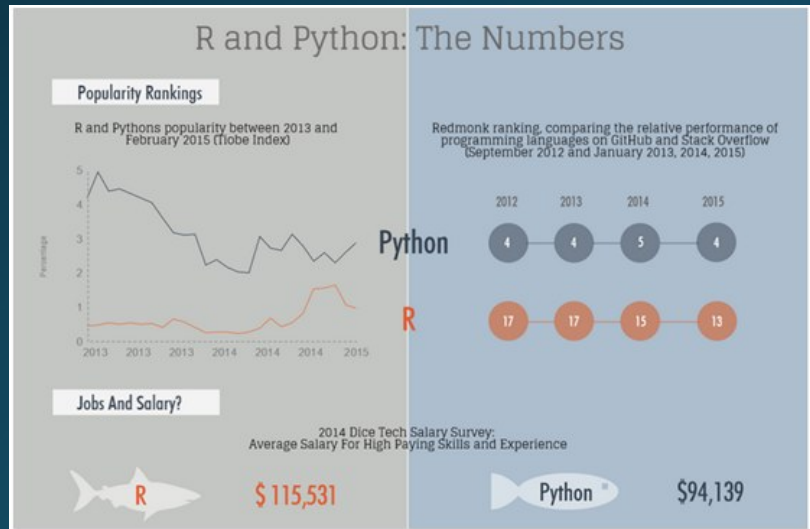


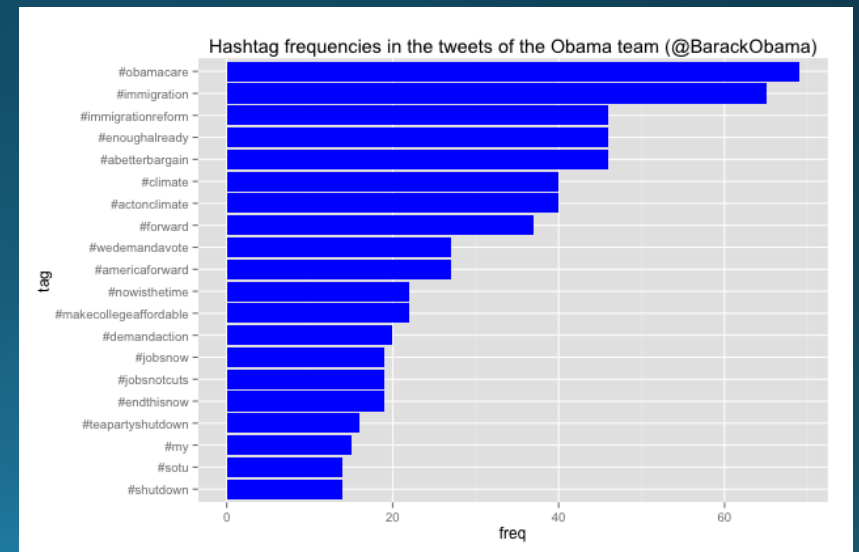
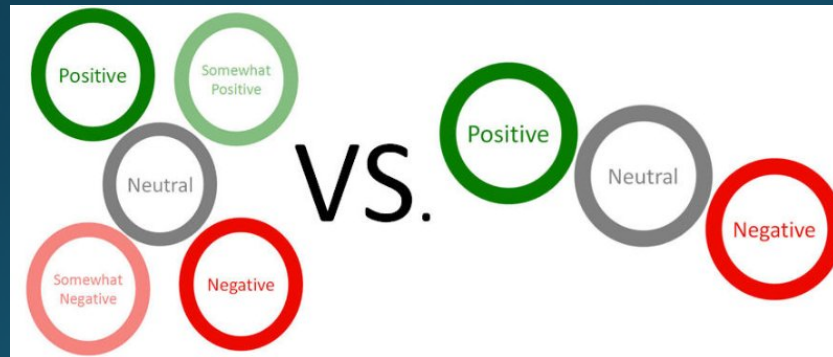
Households with No Motor Vehicle

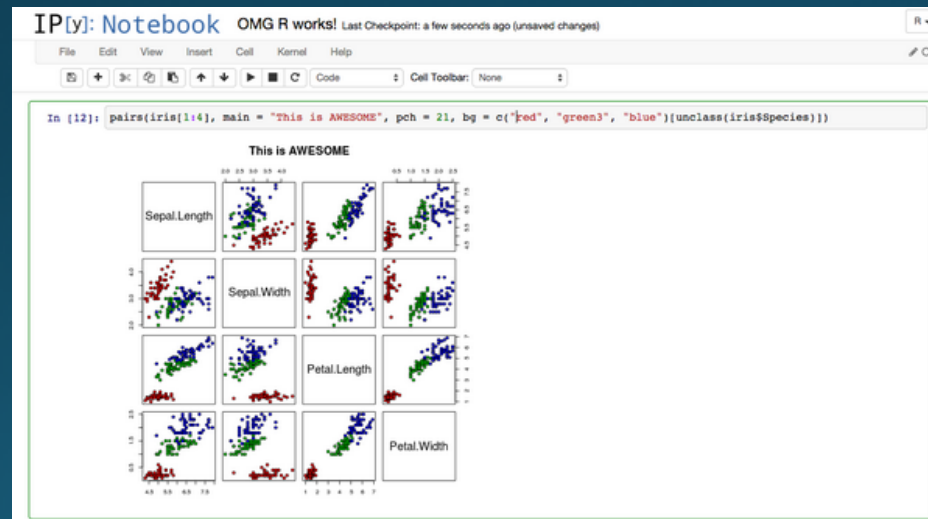
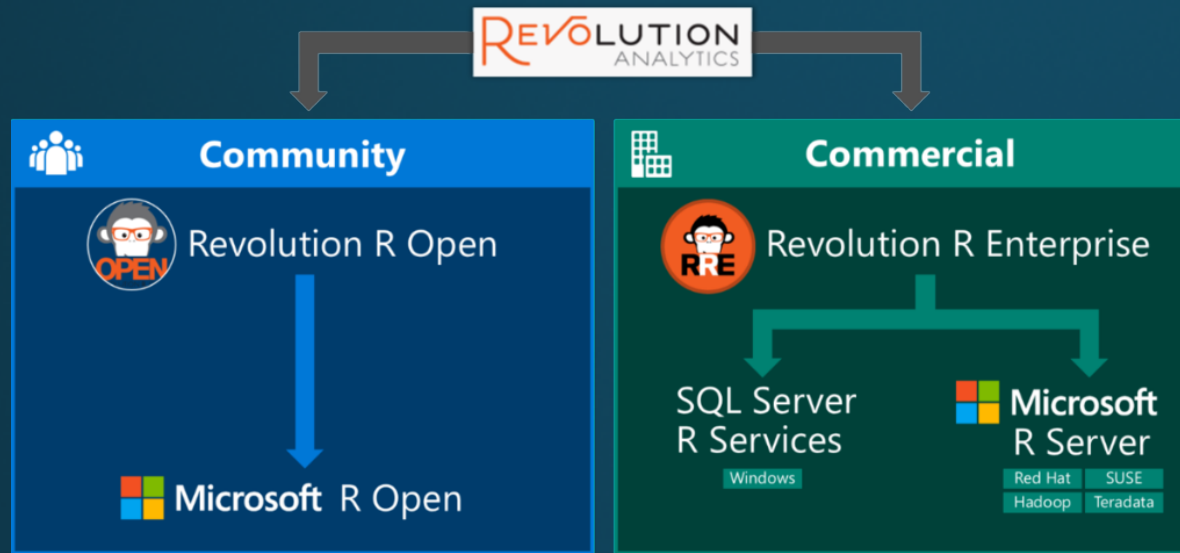
- Orleans Parish reports the highest rate of households with no motor vehicle (18.48%). Orleans Parish includes the City of New Orleans which has more public transportation.

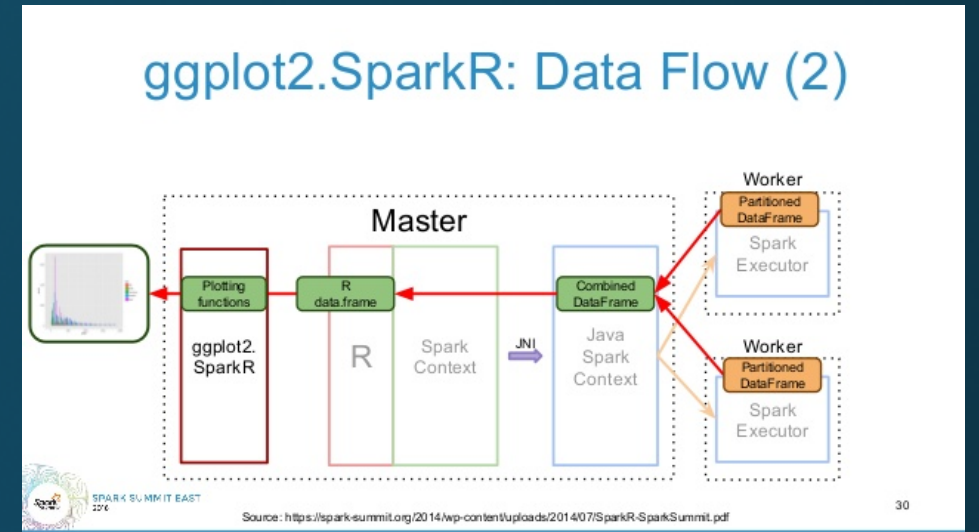


What programming/statistics languages you used for an analytics / data mining / data science work in 2014?	
Language used	<div></div> % voters in 2014 (719 total) <div></div> % voters in 2013 (713 total) <div></div> % voters in 2012 (579 total)
R (352 voters in 2014)	<div></div> 49.0% <div></div> 60.9% <div></div> 52.5%
SAS (262)	<div></div> 36.4% <div></div> 20.8% <div></div> 19.7%
Python (252)	<div></div> 35.0% <div></div> 38.8% <div></div> 36.1%
SQL (220)	<div></div> 30.6% <div></div> 36.6% <div></div> 32.1%
Java (89)	<div></div> 12.4% <div></div> 16.5% <div></div> 21.2%
Unix shell/awk/sed (63)	<div></div> 8.8% <div></div> 11.1% <div></div> 14.7%
Pig Latin/ Hive/ other Hadoop-based languages (61)	<div></div> 8.5% <div></div> 8.0% <div></div> 6.7%
SPSS (58)	<div></div> 8.1% not asked not asked
MATLAB (45)	<div></div> 6.3% <div></div> 12.5% <div></div> 13.1%









Fast!

Statistics!

Scalable

Spark

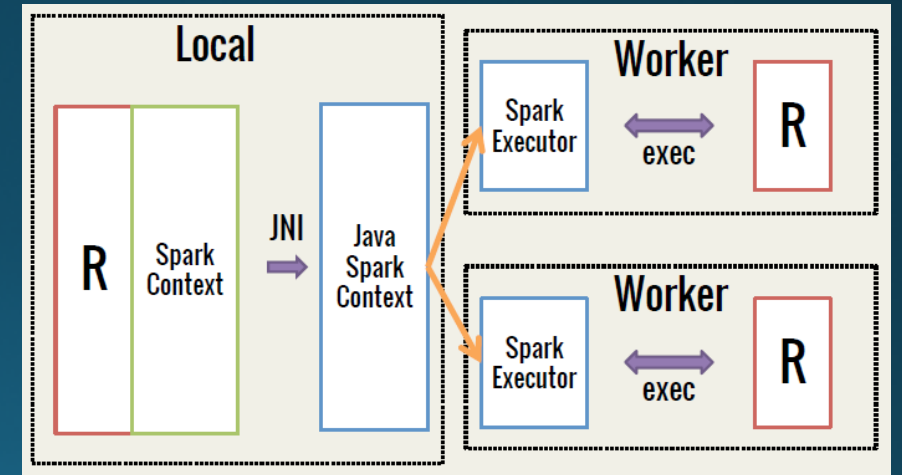
+

R

Packages

Interactive Shell

Plots



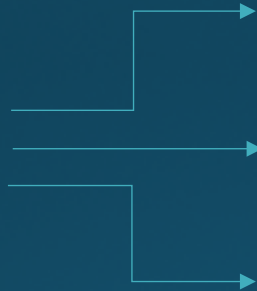
Congratulations!

- You are the newest data scientist at Tesla Motors
- First task: Fuse two data sets
- Skills Required:
- Data cleansing, Missing Value imputation, Classification Modeling, Clustering, Regression Analysis, Data Manipulation

Dataset Link



Dataset A: Tesla
Customers



Dataset B: USA Consumer
Database



That's easy. Use a
common ID!

What if no
common link
exists?

Link Details

Dataset A

ID	Gender M	Gender F	Indivd_Age	Income	Environm entalist
1043	""	female	34	9	1
5021	MALE	""	72	12	1
4863	""	FEMALE	81	6.5	1
4651	male	""	NA	3.2	0
1002	99	99	49	21	0
7523	Male	""	21	-7.9	0

Dataset B

Age	Gender	Income	Environmental
28	M	90	Yes
35	F	81	No
14	F	65	No
65	M	42	Yes
44	M	110	No
39	M	49	No

Learn More

- R Cookbook - Many helpful examples and topics. The “Graphs” section helps tremendously with ggplot2 <http://www.cookbook-r.com/>
- Concise Guide to R – Covers all major R functionality <http://www.cis.upenn.edu/~matuszek/Concise%20Guides/Concise%20R.html>
- Stack Overflow – The best help is available when you post a question here. <http://stackoverflow.com/questions/tagged/r>
- Cross Validated – For statistics and modeling help. <http://stats.stackexchange.com/>
- Introduction to Statistical Learning – The BEST way to learn statistics and machine learning. And it's free!
 - Book: <http://www-bcf.usc.edu/~gareth/ISL/>
 - Video Lectures: <http://www.r-bloggers.com/in-depth-introduction-to-machine-learning-in-15-hours-of-expert-videos/>