# Next steps

Introduction to Data Science with R
www.therbootcamp.com
@therbootcamp

October 2018

# Hello Data Scientist!

In 2 days, 6 sessions, and 16 hours you have come a long way.

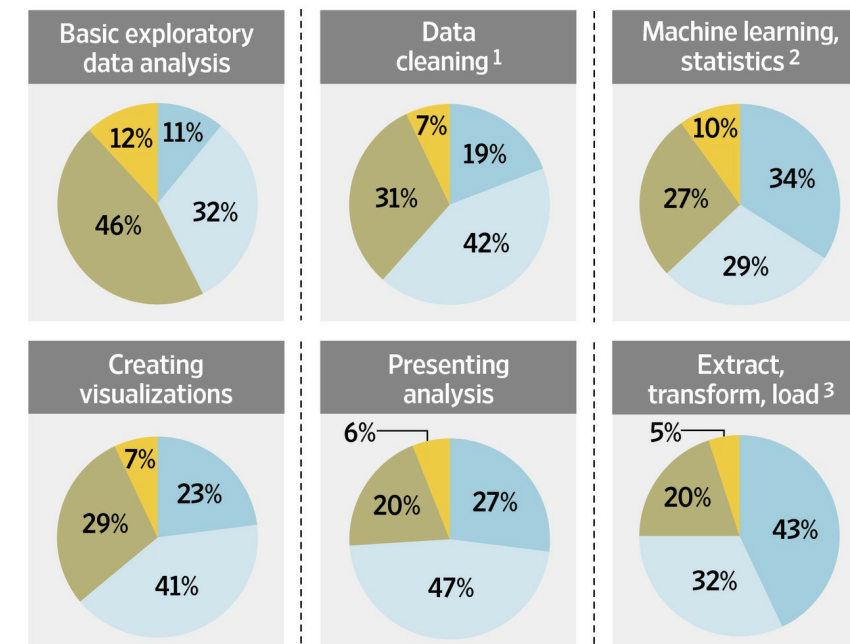| | Sat, 6 Oct | Sun, 7 Oct |
|---|---|---|
| 0900 | Welcome | Recap |
| 0930 | **Intro to R** +Interactive | **Analysing** +Practical |
| 1200 | *Lunch* | *Lunch* |
| 1300 | **Data** +Practical | **Plotting** +Practical |
| 1530 | **Wrangling** +Practical | **Case studies** +Practical |
| 1630 | | Next steps |
| 1800 | Wrapup | *Apero* |

# Data Scientist

Harvard Business Review



**Where Does the Time Go?**

The amount of time spent on various tasks by surveyed nonmanagers in data-science positions

- Less than 1 hour a week
- 1 to 4 hours a week
- 1 to 3 hours a day
- 4 or more hours a day

**Basic exploratory data analysis**
12% 11% 46% 32%

**Data cleaning [1]**
7% 19% 31% 42%

**Machine learning, statistics [2]**
10% 34% 27% 29%

**Creating visualizations**
7% 23% 29% 41%

**Presenting analysis**
6% 20% 27% 47%

**Extract, transform, load [3]**
5% 20% 43% 32%

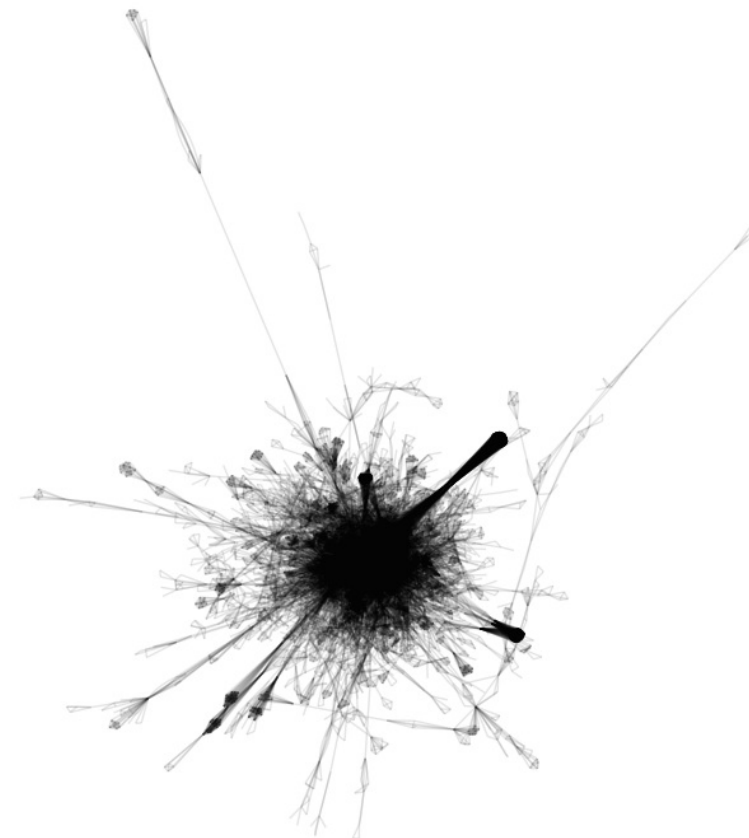[1] Correcting or removing faulty data   [2] Creating computer models
[3] Also known as ETL — moving information to a data warehouse
Source: O'Reilly Media Inc. online survey of more than 600 datascience professionals, conducted from November 2014 to July 2015   THE WALL STREET JOURNAL.

Wall Street Journal

# Next steps

1. Networks
2. Statistics
3. Machine learning
4. Text analysis
5. Rcpp
6. Forms
7. R on servers

8. Books
9. Websites R
10. Groups

source R-bloggers.com

# Stats packages

> "It's easy to lie with statistics; it is easier to lie without them."
>
> Frederick Mosteller

| Package | Description |
|---------|-------------|
| stats | Linear, generalized linear models, individual tests, and distributions. |
| lme4, afex | Mixed-mode, hierarchical regression. |
| sem, lavaan, OpenMx | Structural equation modeling. |
| survival | Survival analysis. |



xkcd.com

# Bayesian statistics

> The subjectivist (i.e. Bayesian) states his judgements, whereas the objectivist sweeps them under the carpet by calling assumptions knowledge, and he basks in the glorious objectivity of science.
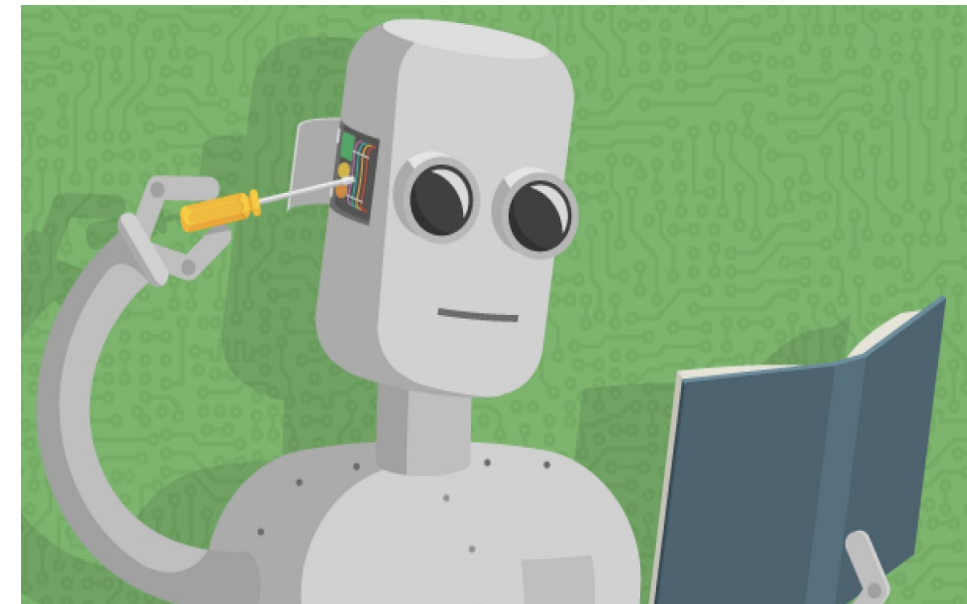>
> I. J. Good

| Package | Description |
|---|---|
| BayesFactor, rstanarm | Bayesian linear models. As easy as non-Bayesian methods. |
| rjags, rstan | Build flexible, hierarchical Bayesian models. |
| mcmc | Metropolis algorithms. |
| bridgesampling | Estimating marginal likelihoods using bridgesampling. |



xkcd.com

# Machine learning

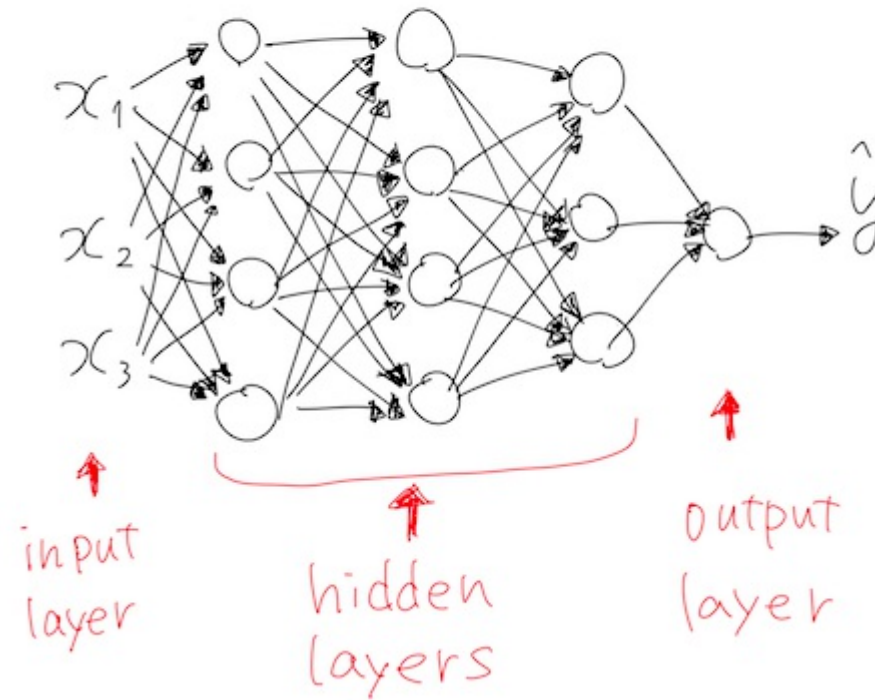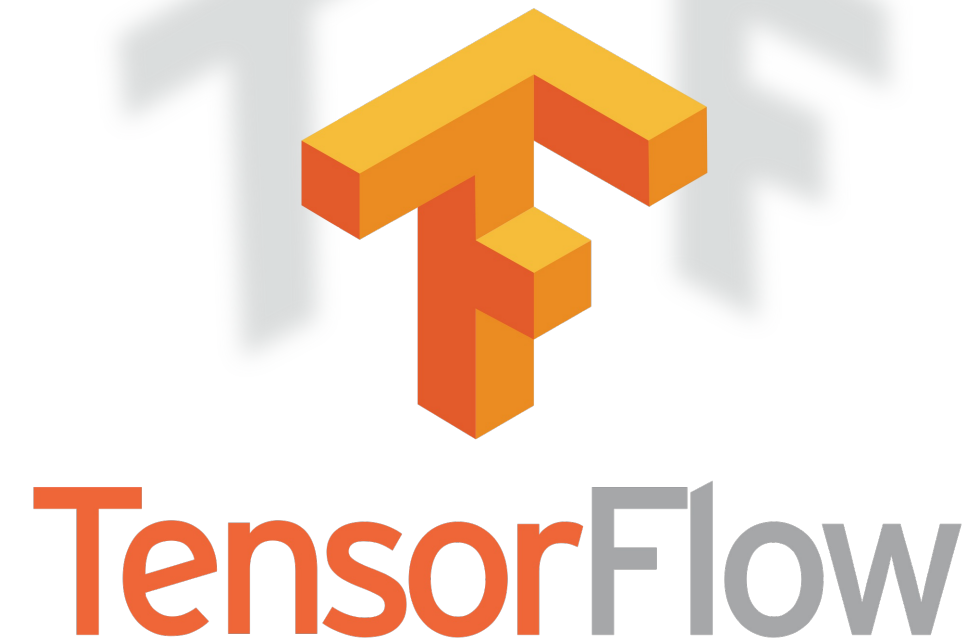| Package | Description |
|---|---|
| caret | Umbrella package for diverse machine learning algorithms. |
| mlr, e1071, etc. | Other umbrella packages. |
| randomForest, rpart, FFTrees | Decision trees. |
| cluster, fastcluster, cstab, etc. | Cluster analysis. |
| forecast, mgm, timeSeries, etc. | Time series models. |
| tensorflow | Time series models. |

# Deep learning

# Deep learning: Tensorflow

```r
# Get Keras
install.packages(keras)
library(keras)
install_keras()

# Define model
model <- keras_model_sequential()
model %>%
  layer_dense(units = 256,
              activation = 'relu',
              input_shape = c(784)) %>%
  layer_dropout(rate = 0.4) %>%
  layer_dense(units = 128,
              activation = 'relu') %>%
  layer_dropout(rate = 0.3) %>%
  layer_dense(units = 10,
              activation = 'softmax')
```
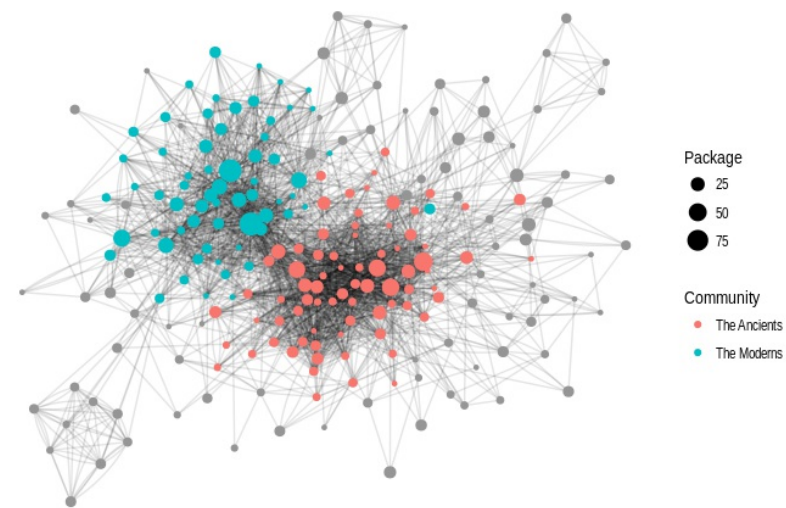
# Text analysis



**Sundar Pichai @ Google IO, May 2018**

| Package | Description |
|---|---|
| tm, tidytext | General text analysis packages |
| stringr, stringi | String operations and regular expressions. |
| rvest, XML | Scraping content of the internet |
| text2vec | Vector representation of words. |
| SentimentAnalysis | Sentiment analysis. |
| twitteR, streamR, jsonlite | Streaming and parsing tweets. |
| Rfacebook | Access to Facebook API. |

# Network analysis

A **social graph** of package Co-authors using `tidyverse` plus `ggraph`, an extension for `ggplot2` for graphs (aka networks) and `igraph`, an extremely powerful library for network analysis. Find the code and additional explanations **here**.



source R-bloggers.com



source R-bloggers.com

# Rcpp

By now one of the most referenced R packages is Rcpp - **R's interface to C++**. With often relatively little effort due to **Rcpp sugar**, Rcpp can provide vast speed improvements, which many packages today rely on Rcpp in the background for **swift code execution**. Rcpp becomes particularly powerful, when supplemented with BH, which makes avaialble a collection **free, peer-reviewed C++ libraries**, and **RcppArmdillo**, which available the high-performance **Armadillo** library for linear algebra methods.



source classic105.com



Quick Reference Guide

# Google Forms & Maps

New packages also allow you to interact with **Google Maps** and **Google Forms**. Use *ggmap* to access Google Maps and *googlesheets* to access Google Forms.

```
library(ggmap)
ggmap(get_map(c(7.588576, 47.559601),zoom=16))
```



## Basel R Bootcamp - Follow-up questionnaire

Please be so kind to take a few minutes and provide us with feedback for the Introduction to Data Science with R bootcamp in October 2018. In the first part you will have the chance to tell us what you think about the workshop in general. In the second part you can indicate which sessions you liked/disliked and why.

*Erforderlich

How did you hear about the bootcamp? *

◯ Google / Google Ad

◯ LinkedIn

◯ Advanced Studies website

◯ Friends and colleagues

◯ Facebook

◯ Sonstiges: _____

# R on servers

With the need for more computing power for, e.g., machine learning, R is also moving more and more to the **server and cluster sphere**. Rstudio offers a convenient (and free) solution to access **RStudio** sessions on a **remote linux server** via a browser. RStdudio also offers the `sparklyr` package integrating the **Spark** architecture for efficient cluster computing with `dplyr`. Other packages exist to interface with **Hadoop** (`RHadoop`) and **SQL** (`RSQLite`).

# Remote databases

R provides **all necessary tools to pull data from or directly work with** remote databases such as, e.g., a SQL database. Find out more at:

db.rstudio.com

# How to continue

# Books

Here is a very incomplete series of good books. They are ordered by complexity, beginning with user-friendly books on **learning statistics** in R and ending with books focusing on the more **advanced topics of the R language**.

# Websites

The web is a great place to learn about R.

**Google** or **Rseek**, which is a wrapper around google to maximize hits related to R. However, most of the time Google works just fine. Just be sure to add   to the the search query.

**R-bloggers** is a website on which R users inform each other on the newest developments. See, e.g., Nathaniel's **entry**.

**Stackoverflow** is a website on which R users exchange problems and solutions to problems. Try post something yourself. You will be amazed by the turnaround.

# R meetup groups

Great place to meet other **R enthusiast** in and outside of industry and see interesting presentions on R's latest developments.

# Consulting

Feel free to contact us.



**Dr. Dirk Wulff**

dirkwulff.org
github.com/dwulff
cstab,
mousetrap, choicepp

**Dr. Nathaniel Phillips**

nathanieldphillips.com
github.com/ndphillips
yarrr,
FFTrees
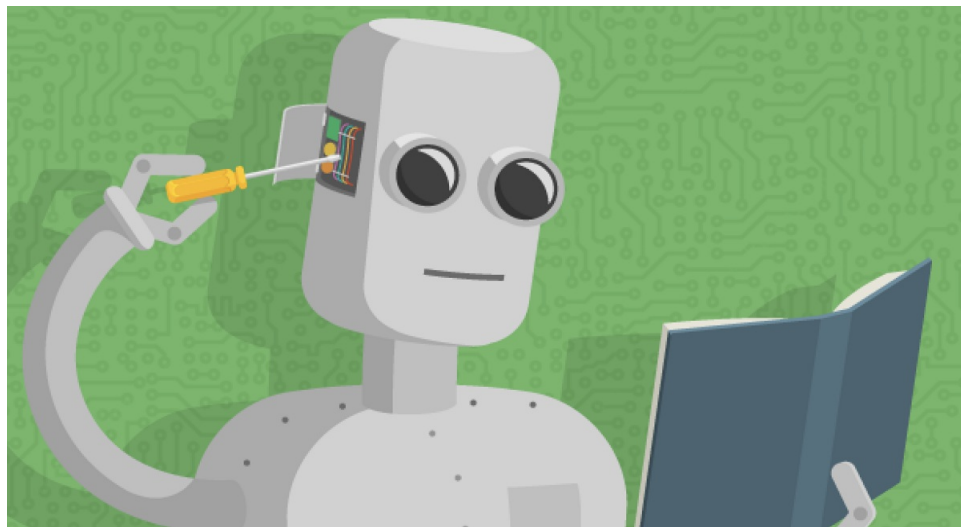
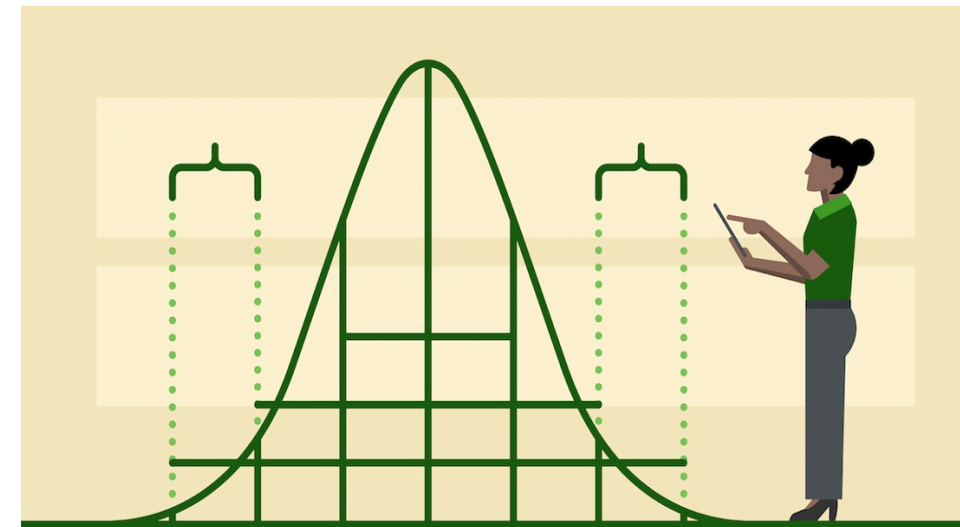**Markus Steiner**

github.com/mdsteiner
ShinyPsych

# Next courses

## Applied Machine Learning with R



**2 days, January 19-20, 2019**

## Statistics with R



**2 days, February/March, 2019**

# Thank you

**Here is an R Joke.**

And now one more thing...