# From laptop to cloud: running R on AWS

Ksenia Aleksankina



travelnest

New data is logged constantly

Model performance needs to be monitored

| Obtaining raw data | → | Cleaning data | → | Exploratory analysis | → | Fitting a model | → | Presenting results |

Data sources and data shape change

Reports need to be updated daily/hourly
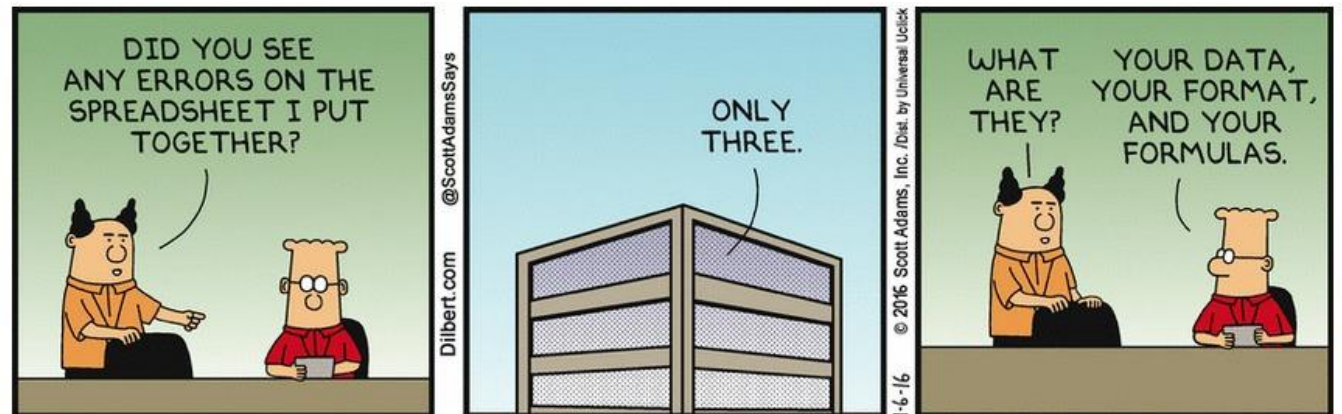
Some of the resulting issues:

- Difficulties in consuming data
- Reproducibility
- Difficulties in communicating analysis results

Solution:

- Deploy and access R/RStudio/Shiny on a cloud server
- Collect, store, and pre-process data in the cloud
- Keep analysis results or model predictions up to date and accessible
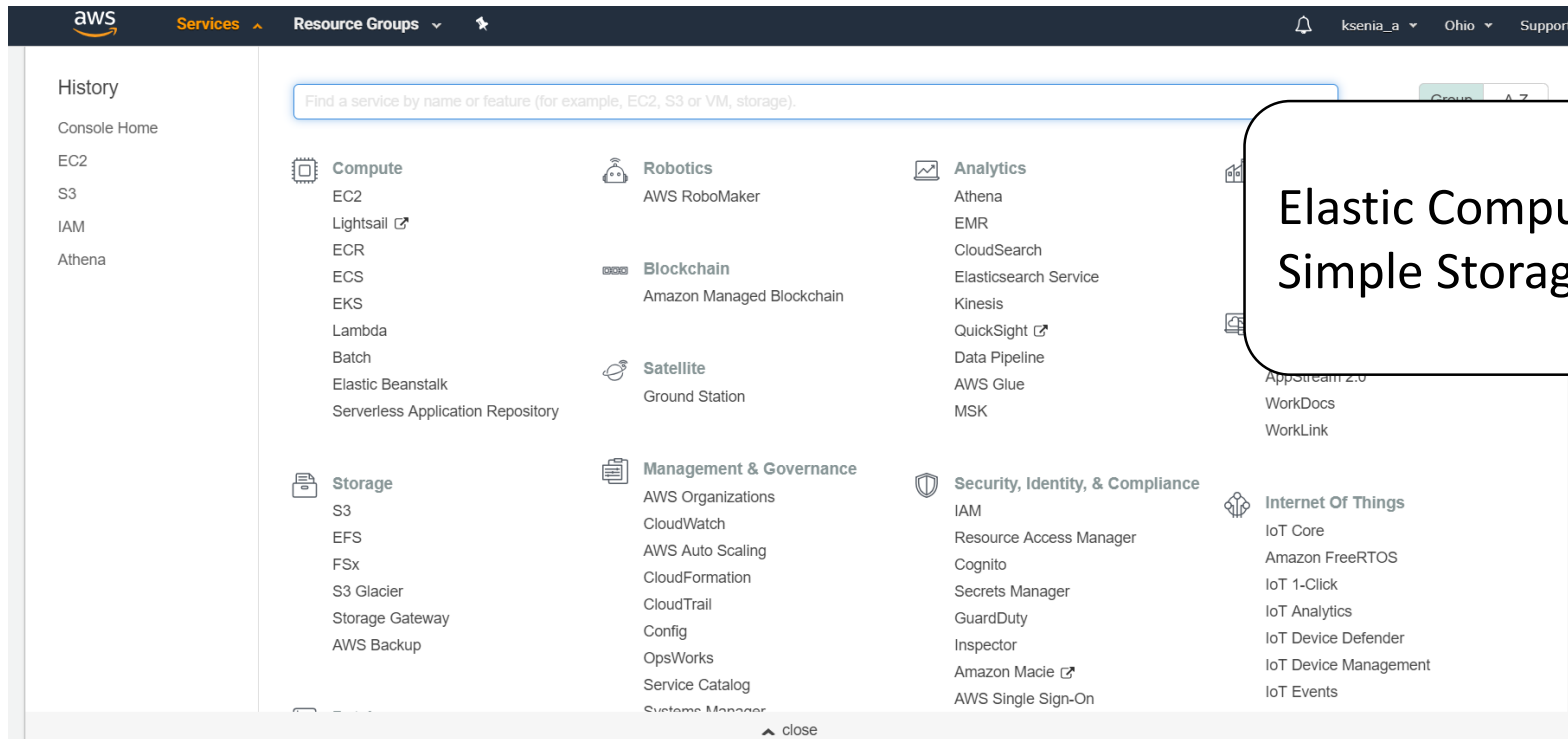
# Why AWS?

Amazon Web Services is an on-demand cloud computing platform

- Variety of building blocks and structures
- Pay-as-you-go (with a free tier)

# The account is created, now what?



Elastic Compute Cloud (EC2)
Simple Storage Service (S3)

https://aws.amazon.com/blogs/big-data/running-r-on-aws/
http://stanke.co/r-on-aws-cloud/

# Before launching EC2 instance

## Create an Identity and Access Management (IAM) role

# Before launching EC2 instance

## Create an Identity and Access Management (IAM) role



AmazonS3FullAccess, AmazonEC2FullAccess

# Before launching EC2 instance

## Create an Identity and Access Management (IAM) role

Create role

Select type of trusted entity

1  2  3

AWS service
EC2, Lambda and others

Anoth
Belong

Allows AWS services to perform actions on your behalf

Choose the service that will use thi

EC2
Allows EC2 instances to call AWS services on your be

Lambda
Allows Lambda functions to call AWS services on you

Create role

▼ Attach permissions policies

Choose one or more policies to attach to your new

Create policy

Filter policies ⌄    🔍 s3

Policy name ▼

☐  ▶  📦 AmazonDMSRedshiftS3Role

☐  ▶  📦 AmazonS3FullAccess

☐  ▶  📦 AmazonS3ReadOnlyAccess

☐  ▶  📦 QuickSightAccessForS3StorageM

AmazonS3FullAccess,

Review and create

Create role

1  2  3  4

Review

Provide the required information below and review this role before you create it.

Role name*  [                    ]

Use alphanumeric and '+=,.@-_' characters. Maximum 64 characters.

Role description  [Allows EC2 instances to call AWS services on your behalf.

                  ]

Maximum 1000 characters. Use alphanumeric and '+=,.@-_' characters.

Trusted entities  AWS service: ec2.amazonaws.com

Policies  📦 AmazonEC2FullAccess ↗
          📦 AmazonS3FullAccess ↗

- Chose an EC2 instance (AMI)
- Instance type (R runs only on one CPU but requires a lot of memory)

- Configure Instance
  Select previously created IAM

| Type | Protocol | Port Range | Source | | Description | |
|---|---|---|---|---|---|---|
| SSH | TCP | 22 | Custom | 0.0.0.0/0 | e.g. SSH for Admin Desktop | ✕ |
| Custom TCP I | TCP | 8787 | Custom | 0.0.0.0/0 | RStudio | ✕ |
| Custom TCP I | TCP | 3838 | Custom | 0.0.0.0/0 | Shiny | ✕ |

Add Rule

- Configure security groups

- 0.0.0.0/0 allow all IP addresses to access your instance

# The instance is ready…



... time to install R, Rstudio, and Shiny

```
# Install R
$ sudo yum install -y R


# Install RStudio Server
$ wget https://download2.rstudio.org/server/centos6/x86_64/rstudio-server-rhel-1.2.1335-x86_64.rpm
$ sudo yum install -y --nogpgcheck /rstudio-server-rhel-1.2.1335-x86_64.rpm


# Install Shiny and Shiny-server
$ R -e \"install.packages('shiny', repos='https://cran.rstudio.com/')\"
$ wget https://download3.rstudio.org/centos6.3/x86_64/shiny-server-1.5.9.923-x86_64.rpm
$ yum install -y --nogpgcheck shiny-server-1.5.9.923-x86_64.rpm


# Add user(s)
$ sudo useradd –m *username*
$ sudo passwd *password*
```
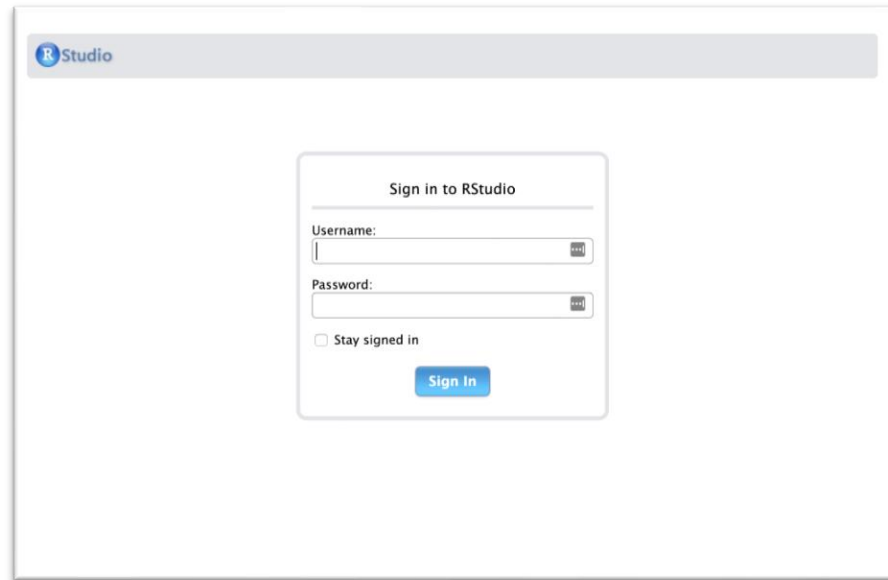
# Issues along the way

- t2.micro did not like dplyr or leaflet

  Solution: change instance type to t2.medium

- Occasional 'non-zero exit status' errors

- For leaflet: Install Portable Network Graphics reference library

  $ sudo yum -y install libpng-devel

# Simple Storage Service (S3)

# RStudio



# Shiny app

# Data for the Shiny app



http://insideairbnb.com/about.html

```r
1    Sys.setenv(
2      "AWS_ACCESS_KEY_ID" = "key",
3      "AWS_SECRET_ACCESS_KEY" = "secret_key",
4      "AWS_DEFAULT_REGION" = "region"
5    )
```

```r
6    library(tidyverse)
7    library(lubridate)
8    library(aws.s3)
9
10   calendar <-
11     read.csv(
12       url(
13         "https://s3.us-east-2.amazonaws.com/ksenia-testbucket/R_test/calendar.csv"
14       )
15     )
16   listings <-
17     read.csv(
18       url(
19         "https://s3.us-east-2.amazonaws.com/ksenia-testbucket/R_test/listings.csv"
20       )
21     )
22   calendar$price = as.numeric(gsub("\\$", "", calendar$price)) * 0.77
23   calendar$date = as.Date(calendar$date)
24
```

```r
27
28   df <-
29     left_join(select(calendar, c(
30       "listing_id", "date", "price", "minimum_nights"
31     )),
32       select(
33         listings,
34         c(
35           "id",
36           "property_type",
37           "accommodates",
38           "neighbourhood_cleansed",
39           "latitude",
40           "longitude"
41         )
42       ),
43       by = c("listing_id" = "id")) %>%
44       filter(accommodates <= 6 & !is.na(price))
45
46   df$month <- round_date(df$date, unit = "month")
47   df$accommodates <-  as.factor(df$accommodates)
48
49   s3write_using(df,
50                 FUN = write.csv,
51                 bucket = "s3://ksenia-testbucket/R_test",
52                 object = "airbnb.csv")
53
```

# Shiny app code

```r
library(shiny)
library(ggplot2)
library(magrittr)
library(lubridate)
library(dplyr)
library(leaflet)

load("./airbnb.Rdata")

neighbourhood <- unique(df[, "neighbourhood_cleansed"])


ui <- pageWithSidebar(
  headerPanel('Airbnb: price per night'),
  sidebarPanel(selectInput('x', 'Neighbourhood', neighbourhood)),
  mainPanel(plotOutput('plot1'),
            leafletOutput('plot2'))
)


server <- function(input, output, session) {
  selectedData <- reactive({
    filter(df, neighbourhood_cleansed == input$x) %>%
      group_by(.dots = c("month", "accommodates")) %>%
      summarise(avg = mean(price),
                n = n(),
                median = median(price))
  })
  selectedData2 <- reactive({
    filter(df,  neighbourhood_cleansed == input$x) %>%
      distinct(listing_id, longitude, latitude)
  })

  output$plot1 <- renderPlot({
    ggplot(selectedData(), aes(month, median, colour = accommodates)) +
      # geom_line(size = 1) +
      geom_point(size = 2.5) +
      labs(title = neighbourhood[1], y = "Median price per night / GBP", x = "Month") +
      theme(text = element_text(size = 20))
  })

  output$plot2 <- renderLeaflet({
    leaflet(selectedData2()) %>%
      addTiles() %>%
      addMarkers( ~ longitude, ~ latitude)
  })
}

shinyApp(ui = ui, server = server)
```

# Useful resources

- Cloudyr project http://cloudyr.github.io/

- Running R on AWS https://aws.amazon.com/blogs/big-data/running-r-on-aws/

- Taking Advanced Analytics to the Cloud http://stanke.co/r-on-aws-cloud/

- What Is Amazon S3? https://docs.aws.amazon.com/AmazonS3/latest/dev/Welcome.html

- An Introduction to Rocker: Docker Containers for R https://arxiv.org/abs/1710.03675

Thank you!