

Big Data Step-by-Step

Boston Predictive Analytics
Big Data Workshop

Microsoft New England Research &
Development Center, Cambridge, MA

Saturday, March 10, 2012

by **Jeffrey Breen**

President and Co-Founder
Atmosphere Research Group
email: jeffrey@atmosgrp.com
Twitter: @JeffreyBreen



<http://atms.gr/bigdata0310>

just need a
little more RAM

~~Big Data~~ Infrastructure

Part 2: Running R + RStudio on Amazon EC2

Code & more on github:

<https://github.com/jeffreybreen/tutorial-201203-big-data>

Overview

- Sometimes you just need a little more RAM, CPU, or disk space than you have
- Let's try launching an instance on Amazon EC2 and configuring it to do some work
- We'll install R and RStudio and call it a day

Some details we'll skip

- Signing up (it's not that hard)

<http://aws.amazon.com/ec2/>

- Pricing (it keeps dropping)

<http://aws.amazon.com/ec2/pricing/>

- The alphabet soup of services (we care about EC2 computing and S3 storage)

Just look for biggest button on the page...

The screenshot shows the AWS Management Console interface for the Amazon EC2 service in the US East (Virginia) region. The main content area is titled 'Amazon EC2 Console Dashboard' and features a 'Getting Started' section with a large yellow box containing the text: 'To start using Amazon EC2 you will want to launch a virtual server, known as an Amazon EC2 instance.' Below this text is a large, prominent blue button labeled 'Launch Instance'. A note below the button states: 'Note: Your instances will launch in the US East (Virginia) region.'

The left sidebar contains a 'Navigation' menu with the following items:

- EC2 Dashboard
- Scheduled Events
- INSTANCES
 - Instances
 - Spot Requests
 - Reserved Instances
- IMAGES
 - AMIs
 - Bundle Tasks
- ELASTIC BLOCK STORE
 - Volumes
 - Snapshots
- NETWORK & SECURITY
 - Security Groups
 - Elastic IPs
 - Placement Groups
 - Load Balancers
 - Key Pairs
 - Network Interfaces

The right sidebar contains a 'My Resources' section with the following statistics:

- 2 Running Instances
- 2 Elastic IPs
- 4 EBS Volumes
- 0 EBS Snapshots
- 2 Key Pairs
- 0 Load Balancers
- 0 Placement Groups
- 6 Security Groups

The bottom of the console shows a 'Service Health' section with a table of service status:

Current Status	Details
✓ Amazon EC2 (US East - N. Virginia)	Service is operating normally

Below the service health table is a link to 'View complete service health details'.

The footer of the console displays the copyright notice: '© 2008 - 2012, Amazon Web Services LLC or its affiliates. All rights reserved.' and links to 'Feedback', 'Support', 'Privacy Policy', 'Terms of Use', and 'An amazon.com company'.

Select an Amazon Machine Image

ami-7385461a is a good, recent 64-bit CentOS image published by RightScale

Request Instances Wizard Cancel

CHOOSE AN AMI

INSTANCE DETAILS

CREATE KEY PAIR

CONFIGURE FIREWALL

REVIEW

Choose an Amazon Machine Image (AMI) from one of the tabbed lists below by clicking its **Select** button.

Quick Start




My AMIs

Community AMIs

Viewing: All Images

ami-7385461a

1 to 1 of 1 Items

AMI ID	Root De	Manifest	Platform	
 ami-7385461a	ebs	411009282317/RightImage_CentOS_5.6_x64_v5.7.14_EBS	 Cent OS	Select 

Only use EBS images

- Instance-storage machines lose their data upon shutdown (termination)
- EBS instances can be stopped and restarted, or terminated when you're done forever

EC2 Amazon Machine Image: ami-7385461a

Description

Tags

AMI ID: ami-7385461a			
AMI Name: RightImage_CentOS_5.6_x64_v5.7.14_EBS			
Description: RightImage_CentOS_5.6_x64_v5.7.14_EBS			
Source: 411009282317/RightImage_CentOS_5.6_x64_v5.7.14_EBS			
Owner: 411009282317	Visibility: Public	Product Code:	
State: available	Kernel ID: aki-825ea7eb	RAM Disk ID:	-
Image Type: machine	Architecture: x86_64	Platform:	Cent OS
Root Device Type: ebs	Root Device: /dev/sda1	Image Size:	8 GiB
Block Devices: /dev/sda1=snap-0eb8296e:8:true			
Virtualization: paravirtual			
State Reason: -			

Pick a size

See <http://aws.amazon.com/ec2/instance-types/>

Request Instances Wizard Cancel

CHOOSE AN AMI **INSTANCE DETAILS** CREATE KEY PAIR CONFIGURE FIREWALL REVIEW

Provide the details for your instance(s). You may also decide whether you want to launch your instances as "on-demand" or "spot" instances.

Number of Instances: **Instance Type:**

☒ **Launch I**
EC2 Instances commonly launched into a VPC.
Launch into:

☐ **Request**

Type	CPU Units	CPU Cores	Memory
Micro (t1.micro)	Up to 2 ECUs	1 Core	613 MB
Large (m1.large)	4 ECUs	2 Cores	7.5 GB
Extra Large (m1.xlarge)	8 ECUs	4 Cores	15 GB
High-Memory Extra Large (m2.xlarge)	6.5 ECUs	2 Cores	17.1 GB
High-Memory Double Extra Large (m2.2xlarge)	13 ECUs	4 Cores	34.2 GB
High-Memory Quadruple Extra Large (m2.4xlarge)	26 ECUs	8 Cores	68.4 GB
High-CPU Extra Large (c1.xlarge)	20 ECUs	8 Cores	7 GB

< Back Continue

Already out of date! Amazon introduced new "m1.medium" instance type this week.

Avoid Premature Termination

Set Termination Protection + Shutdown Behavior

The screenshot shows the 'Request Instances Wizard' in the AWS Management Console. The wizard is at the 'INSTANCE DETAILS' step. The 'Number of Instances' is set to 1 and the 'Availability Zone' is 'No Preference'. Under 'Advanced Instance Options', the 'Kernel ID' and 'RAM Disk ID' are both set to 'Use Default'. The 'Monitoring' checkbox is unchecked. The 'User Data' is set to 'as text'. The 'Termination Protection' checkbox is checked and circled in orange. The 'Shutdown Behavior' dropdown is set to 'Stop' and is also circled in orange. The 'Continue' button is at the bottom right.

Request Instances Wizard Cancel

CHOOSE AN AMI **INSTANCE DETAILS** CREATE KEY PAIR CONFIGURE FIREWALL REVIEW

Number of Instances: 1
Availability Zone: No Preference

Advanced Instance Options
Here you can choose a specific [kernel](#) or [RAM disk](#) to use with your instances. You can also choose to enable CloudWatch Detailed Monitoring or enter data that will be available from your instances once they launch.

Kernel ID: Use Default **RAM Disk ID:** Use Default

Monitoring: ☐ Enable CloudWatch detailed monitoring for this instance
(additional charges will apply)

User Data:
☒ as text
☐ as file ☐ base64 encoded

Termination Protection: ☒ Prevention against accidental termination.

Shutdown Behavior: Stop
Stop
Terminate
avior when the instance is within the instance.

[< Back](#) [Continue](#)

Name your instance

Request Instances WizardCancel

✓

○

CHOOSE AN AMI **INSTANCE DETAILS** CREATE KEY PAIR CONFIGURE FIREWALL REVIEW

Add tags to your instance to simplify the administration of your EC2 infrastructure. A form of metadata, tags consist of a case-sensitive key/value pair, are stored in the cloud and are private to your account. You can create user-friendly names that help you organize, search, and browse your resources. For example, you could define a tag with key = Name and value = Webserver. You can add up to 10 unique keys to each instance along with an optional value for each key. For more information, go to [Using Tags](#) in the *EC2 User Guide*.

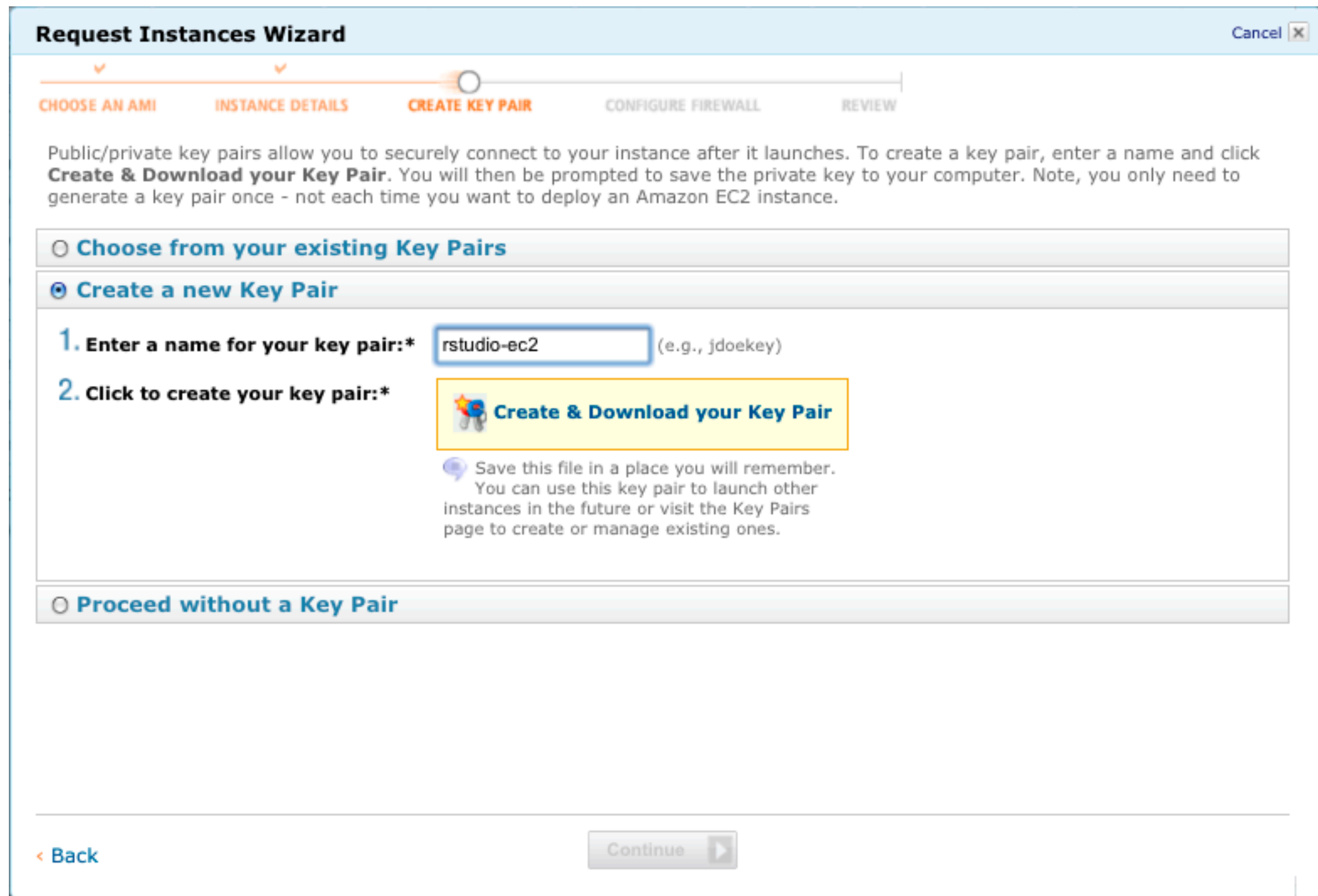
Key (127 characters maximum)	Value (255 characters maximum)	Remove
Name	rstudio	✗
		✗

[Add another Tag.](#) (Maximum of 10)

< BackContinue >

Create a key pair

Don't forget to download it (and keep it safe!)



Request Instances Wizard Cancel

✓ ✓ ○ | CONFIGURE FIREWALL | REVIEW


CHOOSE AN AMI | INSTANCE DETAILS | **CREATE KEY PAIR**


Public/private key pairs allow you to securely connect to your instance after it launches. To create a key pair, enter a name and click **Create & Download your Key Pair**. You will then be prompted to save the private key to your computer. Note, you only need to generate a key pair once - not each time you want to deploy an Amazon EC2 instance.

☐ Choose from your existing Key Pairs

☒ **Create a new Key Pair**

1. Enter a name for your key pair:* (e.g., jdoekey)

2. Click to create your key pair:*  **Create & Download your Key Pair**

 Save this file in a place you will remember. You can use this key pair to launch other instances in the future or visit the Key Pairs page to create or manage existing ones.

☐ Proceed without a Key Pair

[< Back](#) Continue

Create a Security Group

All TCP, UDP and ICMP from your IP address

Request Instances Wizard

Cancel

CHOOSE AN AMI

INSTANCE DETAILS

CREATE KEY PAIR

CONFIGURE FIREWALL

REVIEW

Security groups determine whether a network port is open or blocked on your instances. You may use an existing security group, or we can help you create a new security group to allow access to your instances using the suggested ports below. Add additional ports now or update your security group anytime using the Security Groups page.

☐ Choose one or more of your existing Security Groups

☒ Create a new Security Group

Group Name

rstudio

Group Description

only from my IP

Inbound Rules

Create a new rule:

All ICMP

Source:

74.104.166.41/32

(e.g., 192.168.2.0/24, sg-47ad482e, or 1234567890/default)

+

 Add Rule

TCP

Port (Service)	Source	Action
0 - 65535	74.104.166.41/32	Delete

UDP

Port (Service)	Source	Action
0 - 65535	74.104.166.41/32	Delete

< Back

Continue

Saturday, March 10, 2012

Don't know your IP address?

Don't ask me. Ask Google!



(simply append “/32” when entering into firewall rules)

3... 2... 1...

Request Instances WizardCancel

CHOOSE AN AMI


INSTANCE DETAILS

CREATE KEY PAIR

CONFIGURE FIREWALL

REVIEW

Please review the information below, then click **Launch**.


AMI:  Cent OS AMI ID ami-7385461a (x86_64) [Edit AMI](#)

Number of Instances: 1
Availability Zone: No Preference
Instance Type: Large (m1.large)
Instance Class: On Demand [Edit Instance Details](#)

Monitoring: Disabled **Termination Protection:** Enabled
Tenancy: Default
Kernel ID: Use Default **Shutdown Behavior:** Stop
RAM Disk ID: Use Default
User Data: [Edit Advanced Details](#)

Key Pair Name: rstudio-ec2 [Edit Key Pair](#)

Security Group(s): sg-cfe234a7 [Edit Firewall](#)

[< Back](#) [Launch](#) 

State = running

Up and running at specified domain name

My Instances

Launch Instance Instance Actions Show/Hide Refresh Help

Viewing: All Instances All Instance Types Search 1 to 4 of 4 Instances

	Name	Instance	AMI ID	Root Device	Type	State	Status Checks	Monitoring	Security Groups
<input checked="" type="checkbox"/>	rstudio	i-4f59882b	ami-7385461a	ebs	t1.micro	● running	✓ 2/2 checks passed	basic	rstudio

1 EC2 Instance selected.

EC2 Instance: rstudio (i-4f59882b) ec2-107-22-109-130.compute-1.amazonaws.com

Description Status Checks Monitoring Tags

AMI: RightImage_CentOS_5.6_x64_v5.7.14_EBS (ami-7385461a)

Security Groups: rstudio

State: running

Owner: 581302678308

Subnet ID: -

Virtualization: paravirtual

Reservation: r-a9fbc9

Platform: -

Kernel ID: aki-825ea7eb

AMI Launch Index: 0

Root Device: sda1

Tenancy: default

Lifecycle: normal

Block Devices: sda1

Zone: us-east-1b

Type: t1.micro

Scheduled Events: No scheduled events

VPC ID: -

Source/Dest. Check:

Placement Group:

RAM Disk ID: -

Key Pair Name: rstudio-ec2

Monitoring: basic

Elastic IP: -

Root Device Type: ebs

Time to get all command line

- You'll need an ssh client and the key pair we generated in order to connect with your instance
- We'll use the Cloudera VM to control versions, options, etc.
- ssh won't use your key pair if its file permissions are too lax

```
$ chmod og-rwx rstudio-ec2.pem
```

- Log in as root to your domain name

```
$ ssh -i rstudio-ec2.pem root@YOURDOMAINHERE.amazonaws.com  
(from previous slide)
```

Install R and RStudio

- Create a user login for yourself (RStudio needs this)

```
# useradd jbrene
```

```
# passwd jbrene
```

- EPEL is already installed, so R is easy

```
# yum -y install R
```

- Follow RStudio's download instructions

<http://www.rstudio.org/download/server>

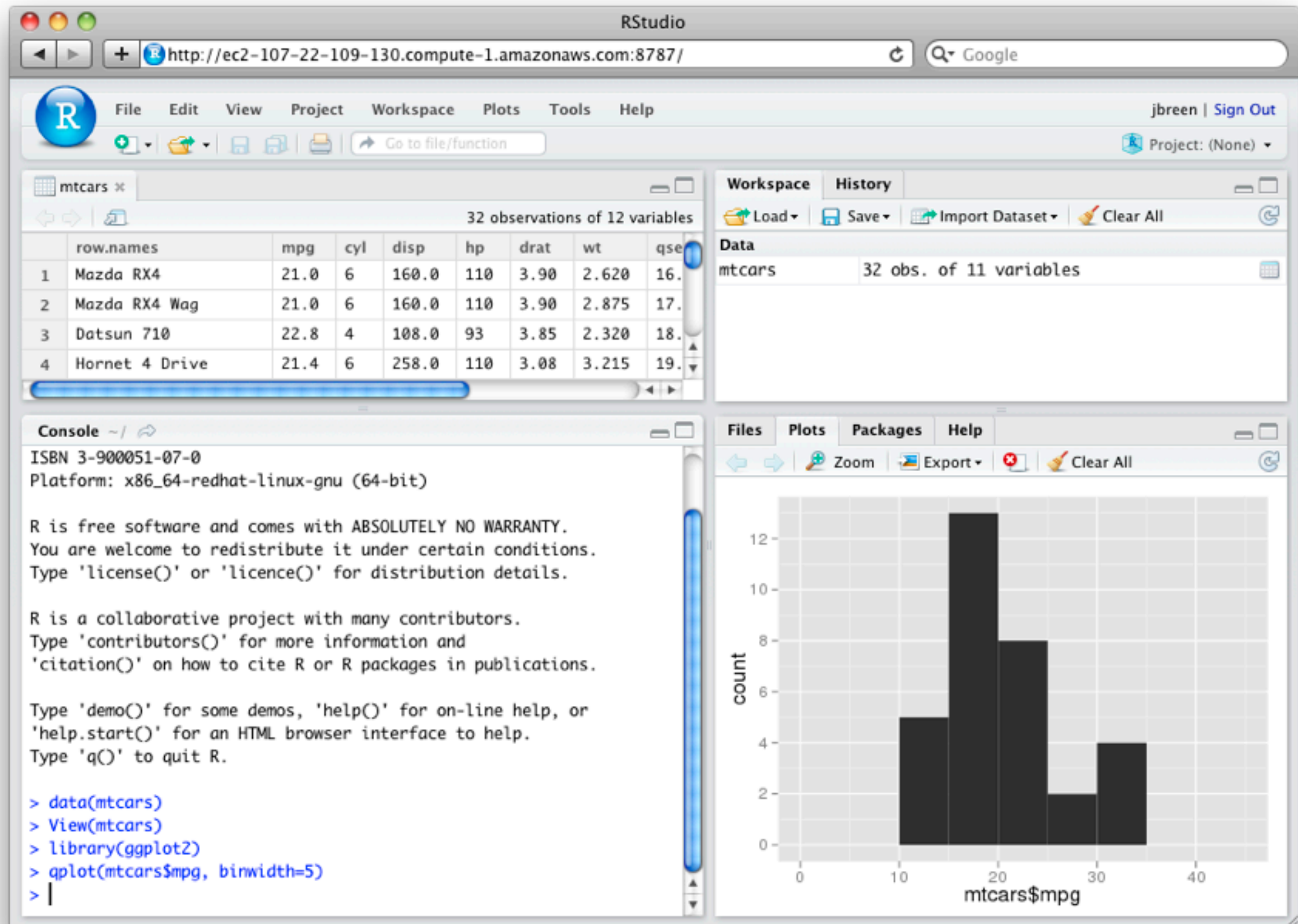
```
# wget http://download2.rstudio.org/rstudio-server-0.95.262-x86_64.rpm
```

```
# rpm -Uvh rstudio-server-0.95.262-x86_64.rpm
```

- Browse to port 8787 and use the login and password

e.g., <http://ec2-107-22-109-130.compute-1.amazonaws.com:8787/>

Success!



The meter's running

- Amazon charges by the hour (or fraction thereof). So when you're done, you should probably shutdown
- via command line

```
$ sudo shutdown -h now
```
- or with the “Stop” Instance Action in the AWS Management Console
- (use “Terminate” if you never want to use it again)

Next up:

How to launch Hadoop
clusters in the cloud
without really trying