

RsparkHub: Scaling RSpark with Kubernetes

Jim Harner, WVU, Mark Lilback, Rc²ai, and Alex Harner, Rc²ai

SDSS 2019

Important links

SDSS 2019 Talk:

<https://github.com/jharner/SDSS2019Docker-k8s>

RSpark: <https://github.com/jharner/rspark>

RSpark Tutorial: <https://github.com/jharner/rspark-tutorial>

Rc² Server: <https://github.com/rc2server>

Rc² Swift Client: <https://github.com/mlilback/rc2SwiftClient>

DevOps for Data Science

Data Science Platforms are complex! How do we solve the data analysis issues of:

- ▶ collaboration,
- ▶ sharing,
- ▶ reliability
- ▶ scalability, and
- ▶ reproducibility?

DevOps (development + operations) has been revolutionized by containers and methods of orchestrating containers. Specifically, we will examine two technologies which are changing how complex software systems are built, deployed, and maintained. These are:

- ▶ Docker: a technology for building, deploying, and running container images;
- ▶ Kubernetes: a technology for deploying, scaling, and managing containerized applications.

Docker Swarm is not as scalable or feature complete as Kubernetes.

Docker

Docker containers are becoming essential for software development and for the deployment of complex computing environments. This talk presents Docker solutions to address these and other use cases relevant to statisticians and data scientists.

- ▶ Rocker: R-based images built by and hosted on Docker Hub
- ▶ RSpark: R and Spark-based images hosted on Docker Hub
- ▶ Rc²: R-based image hosted on GitHub.

Perhaps the most important reason for using Docker by statisticians and data scientists is the complexity of their computing environments, e.g., aligning version numbers for R and its packages, drivers, SQL and NoSQL databases, Hadoop, Spark and its packages, TensorFlow and GPUs, etc.

What is Docker?

Docker allows developers, devops, and sysadmins to develop, deploy, and run applications using containers. We call this containerization.

Containers are:

- ▶ Flexible
- ▶ Lightweight
- ▶ Portable
- ▶ Scalable

A container runs natively on Linux and shares the kernel of the host machine with other containers. A container runs as a discrete process and thus its memory requirements are nearly equivalent to other executables, i.e., it is lightweight. On the other hand, a virtual machine runs a guest OS which is built on a hypervisor, through which host resources are accessed. Running multiple VMs is very heavy.

Why Use Docker?

A Container Platform provides a complete solution, e.g., the components needed for teaching a data science program. Increasingly these platforms are being built on Kubernetes.

Images and Containers

The two principal Docker entities are:

- ▶ Image: an executable package that includes everything needed to run an application
- ▶ Container: a runtime instance of an image

The image contains the code, configuration files, environmental variables, libraries, and the runtime. You can see the images by running the following command in a terminal:

```
docker images
```

A container is an image with state, i.e., a user process. You can view the running containers by executing:

```
docker ps
```

or you can see all containers by executing:

```
docker ps -a
```

Running an Image

You can run an image by:

```
docker run hello-world
```

If the `hello-world` image is not found locally, it is automatically downloaded and launched from Docker Hub, the default image registry.

Rocker Project

Rocker project: a widely-used suite of Docker images with customized R environments for particular tasks.

Currently, there are 17 repos in rocker-org.

The most important to us are:

- ▶ rocker-versioned
 - ▶ r-ver: versioned base R
 - ▶ rstudio: adds rstudio
 - ▶ tidyverse: adds tidyverse and devtools
 - ▶ verse: adds tex and publishing packages
- ▶ shiny
- ▶ ml

These and other available images provide a large number of use cases.

Running R/Rstudio in Containers

A terminal version of R can be run in `r-base` or `r-ver`.

```
docker run --rm -ti rocker/r-base bash
```

The `rocker/r-base` image is downloaded from Docker Hub (if not already in the local filesystem) and run. The `-it` flags cause a `bash` session to start. Typing `R` into the terminal starts `R` (currently version 3.5). The `-rm` flag causes the container to be removed when the session is exited.

Or just:

```
docker run --rm -ti rocker/r-base
```

RStudio can be run by:

```
docker run --rm -p 8787:8787 rocker/rstudio
```

Open a browser with the URL `localhost:8787` and sign in with user/password: `rstudio/rstudio`.

Extending Images

A developer or user may find that an image needs to be extended, e.g., a critical R package is missing or needs to be updated. This can be done in one of two ways:

- ▶ install the R package inside the container using `install.packages`, or
- ▶ rebuild the image by modifying the Dockerfile.

The first method results in mutable infrastructure, i.e., changes are made incrementally. The state of the infrastructure is thus dependent on incremental updates. This method should only be used in critical situations and should be viewed as a temporary fix.

The second method results on an immutable infrastructure, i.e., the state does not change based on actions by the user. Rather, the image is completely rebuilt by the developer.

Rocker Deployment

Any of the Rocker images can be deployed in a single container. This container can be run on your laptop, on your departmental server, or on a cloud platform such as AWS. Although the Rocker images can be built from their GitHub sources, typically the pre-built images on Docker Hub are simply downloaded and run.

Why go to this trouble? Why not just run RStudio locally or on a server?

Suppose you are writing a paper and you want it to be truly reproducible. At this point you want to capture and maintain not only the Rmd or Rnw files, but also the synchronized R and required package versions. Further, the R environment can be synchronized with Linux distro and version along with required storage platform for your data.

But the reasons go well beyond reproducibility.

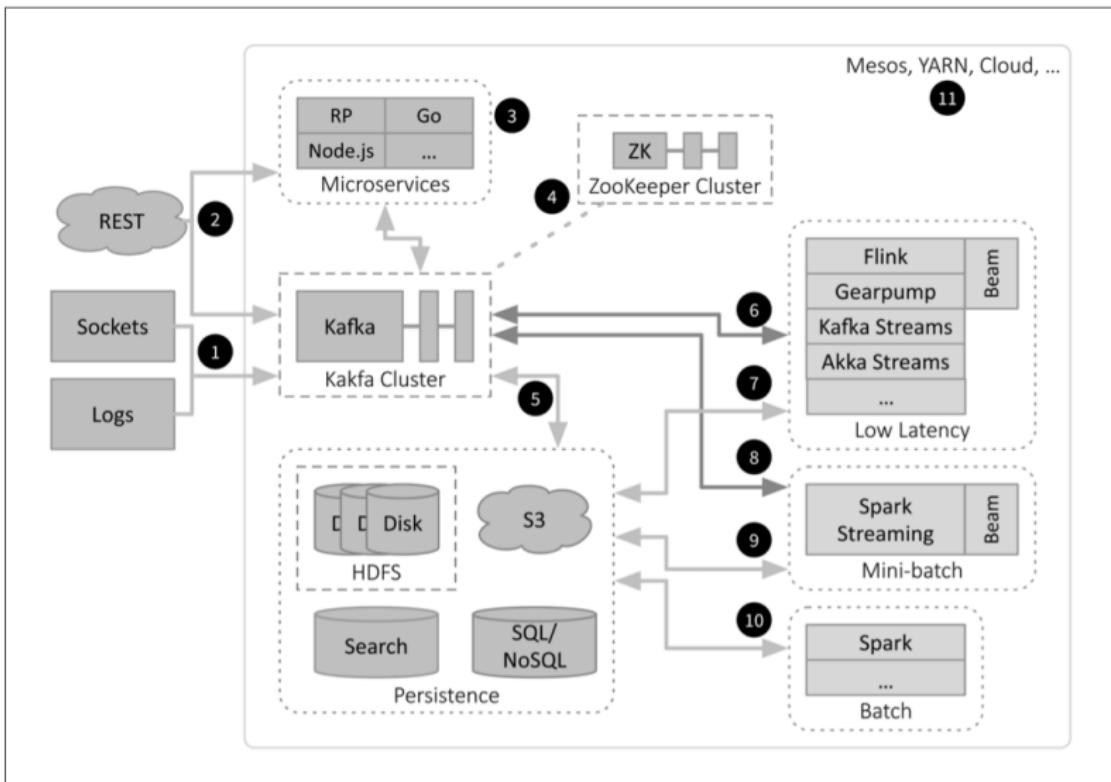
Using Docker for Continuous Integration Deployment [DevOps]

Docker is a platform for containerized networking, compute, and storage for distributed applications. Docker improves:

- ▶ Velocity: increase flow by continuous integration (CI) and continuous delivery (CD);
- ▶ Variation: converge artifacts with immutable binaries;
- ▶ Visualization: visualize (elevate) services in the pipeline by bounding context and using containers for isolation.

Reference paper on CI/CD

A Modern Data Science Platform



Rspark Project

RSpark extends Rocker to embrace big-data platforms. In particular, RSpark's main image, `rstudio` is built on Rocker's `verse` image, which contains RStudio, the Tidyverse, LaTeX, etc. Various data science R packages are added to this image, e.g., `sparklyr`, `SparkR`, `tensorflow`, `rhipe`, etc.

RSpark has other images commonly used in data science:

- ▶ PostgreSQL: a powerful relational database;
- ▶ Hadoop: a distributed storage (HDFS) and computing (MapReduce) environment;
- ▶ Hive: a distributed data warehouse built on HDFS;
- ▶ Spark: a cluster version for distributed computing.

Note: RSpark could be built using RCloud rather than RStudio as an interface to R. RCloud is built on a lower-level image, i.e., `rocker/dr` than the `rocker/verse` image used in RSpark.

We can scale the Spark deployment to multiple workers in RSpark. But how do we scale RSpark to a large number of workers for

Kubernetes

Kubernetes is a platform for running and coordinating related, containerized applications across a cluster of (typically virtual) machines. It manages the complete lifecycles of containerized applications and services to provide:

- ▶ predictability
- ▶ reliability
- ▶ scalability
- ▶ availability

Kubernetes Principles

Container and Kubernetes developers should adhere to certain principles:

- ▶ immutable infrastructure: update an image with a newer image in a single operation, i.e., no incremental updates (except in emergencies);
- ▶ declarative configuration: describe the desired state of your application.

In mutable infrastructures changes are applied incrementally, e.g., `apt-get` in Linux or `install.packages` in R. This should only be done as a temporary fix. An imperative declaration defines the state as a series of actions rather than a declaration of the desired state as is done with Kubernetes.

Kubernetes makes sure the actual state of your environment is aligned with the desired state. Not only that, Kubernetes is an online, self-healing system, i.e., it continuously takes action to ensure the current state matches the desired state.

Kubernetes Architecture

The machines in a Kubernetes cluster are classified as:

- ▶ a master server: the primary point of contact with the cluster. It exposes the Kubernetes API, aligns the actual state to the desired state, schedules work, orchestrates communication, etc.
- ▶ nodes: servers responsible for running workloads using local and external resources.

The API server is the management point of the cluster, e.g., for configuring Kubernetes workloads and organizational units. The API interface is RESTful. A CLI client called `kubectl` is the default method of interacting with the Kubernetes cluster from a local computer.

The basic Kubernetes unit is the pod. Containers are assigned to pods—not hosts. Pods should contain a single application or related applications, e.g., ones that need to share the same filesystem.

RSpark on Kubernetes

RSpark consists of an ensemble of containers, but the `rstudio` container is central and provides RStudio Server and various drivers allowing connectivity to other services, e.g., Spark and HDFS. However, RStudio Server Pro and Launcher are required to run Kubernetes. See: [Using Docker images with RStudio Server Pro, Launcher, and Kubernetes](#).

The intent of RSparkHub is to launch RSpark for each user. However, there appears to be no way to shutdown RStudio processes on demand so a 0 exit code is returned to Kubernetes. Instead, we have automated the process using Ansible scripts on AWS, which spins up an instance for each user.

Rc² on Kubernetes

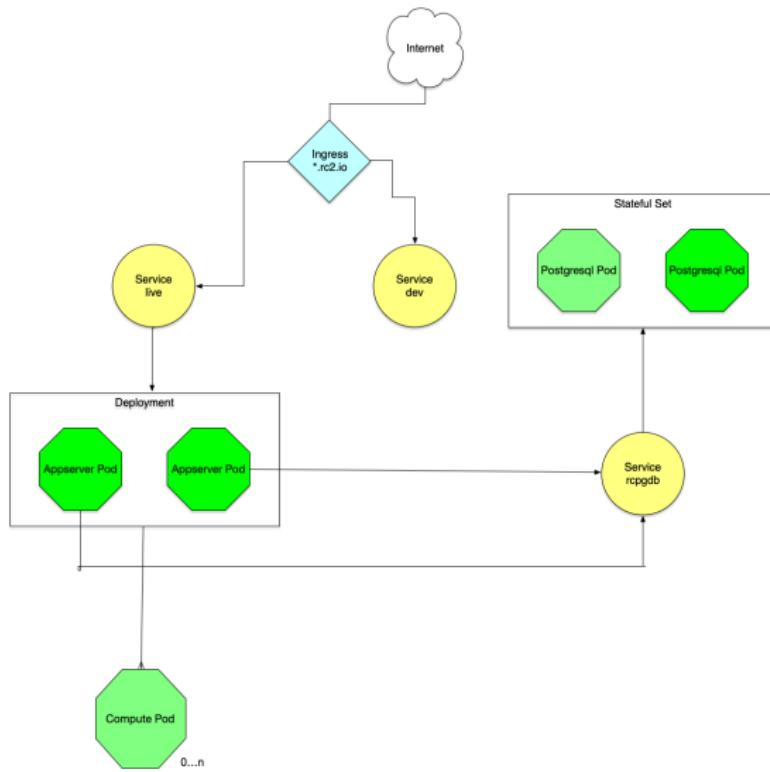
Rc² is a client-server environment for easily scaling R in a Kubernetes cluster. Other languages will be supported, e.g., Python and bash. Various services beyond R can be invoked, e.g., Spark and TensorFlow.

Rc² has two components:

- ▶ Rc² Server: a modern architecture for running R projects designed to interface with other environments at scale, e.g., Spark and Tensorflow
- ▶ Rc² Client: written in Swift currently running on macOS and iOS.

Rc² is being deployed to Digital Ocean Kubernetes.

Rc² Architecture



Rc² Client with an R File

The screenshot shows the Rc² ClientR interface. On the left, a sidebar lists "Source Files" including dframe.R, mdtest.Rmd, newton.Rmd, and puromycin.R. The "puromycin.R" file is currently selected. The main area displays the R code and its execution results.

```
default
Source Files
dframe.R
mdtest.Rmd
newton.Rmd
puromycin.R
Other
mdtest.html
newton.html

1 conc <- c(.02, .02, .06, .06, .11, .11, .22, .22, .56, .56, 1.1, 1.1)
2 velocity <- c(76L, 47L, 97L, 107L, 123L, 139L, 159L, 152L, 191L,
201L, 207L, 208L)
3
4 puromycin.nls <- nls(velocity ~ (beta0 * conc) / (beta1 + conc),
+ start=c(beta0=200, beta1=0.1))
5 summary(puromycin.nls)
6
7 plot(conc, velocity, xlab="Concentration", ylab="Velocity")
8 conc.seq <- seq(min(conc), max(conc), length=100)
9 lines(conc.seq, predict(puromycin.nls, list(conc=conc.seq)), col =
"red", lwd = 2)
10
11 puromycin.df <- data.frame(conc=conc, velocity=velocity)
12 puromycin.df
13 str(puromycin.df)
14
15 ls()
16 rm(conc, velocity)
17 ls()
18
19 puromycin.nls <- nls(velocity ~ (beta0 * conc) / (beta1 + conc),
20                         data=puromycin.df, start=c(beta0=200,
21 beta1=0.1))
22 puromycin.nls$data
23 puromycin.nls
```

The R code defines a non-linear least squares model (nls) for the relationship between concentration and velocity. It includes a plot of the data and a fitted curve. The summary output shows the estimated parameters and their significance.

```
> conc <- c(.02, .02, .06, .06, .11, .11, .22, .22, .56, .56, 1.1, 1.1)
> velocity <- c(76L, 47L, 97L, 107L, 123L, 139L, 159L, 152L, 191L,
201L, 207L, 208L)

> puromycin.nls <- nls(velocity ~ (beta0 * conc) / (beta1 + conc),
+ start=c(beta0=200, beta1=0.1))

> summary(puromycin.nls)

Formula: velocity ~ (beta0 * conc)/(beta1 + conc)

Parameters:
            Estimate Std. Error t value Pr(>|t|)
beta0 2.127e+02 6.947e+00 30.615 3.24e-11 ***
beta1 6.412e-02 8.281e-03 7.743 1.57e-05 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 10.93 on 10 degrees of freedom

Number of iterations to convergence: 6
Achieved convergence tolerance: 6.085e-06

> plot(conc, velocity, xlab="Concentration", ylab="Velocity")
> conc.seq <- seq(min(conc), max(conc), length=100)
> lines(conc.seq, predict(puromycin.nls, list(conc=conc.seq)), col =
"red", lwd = 2)

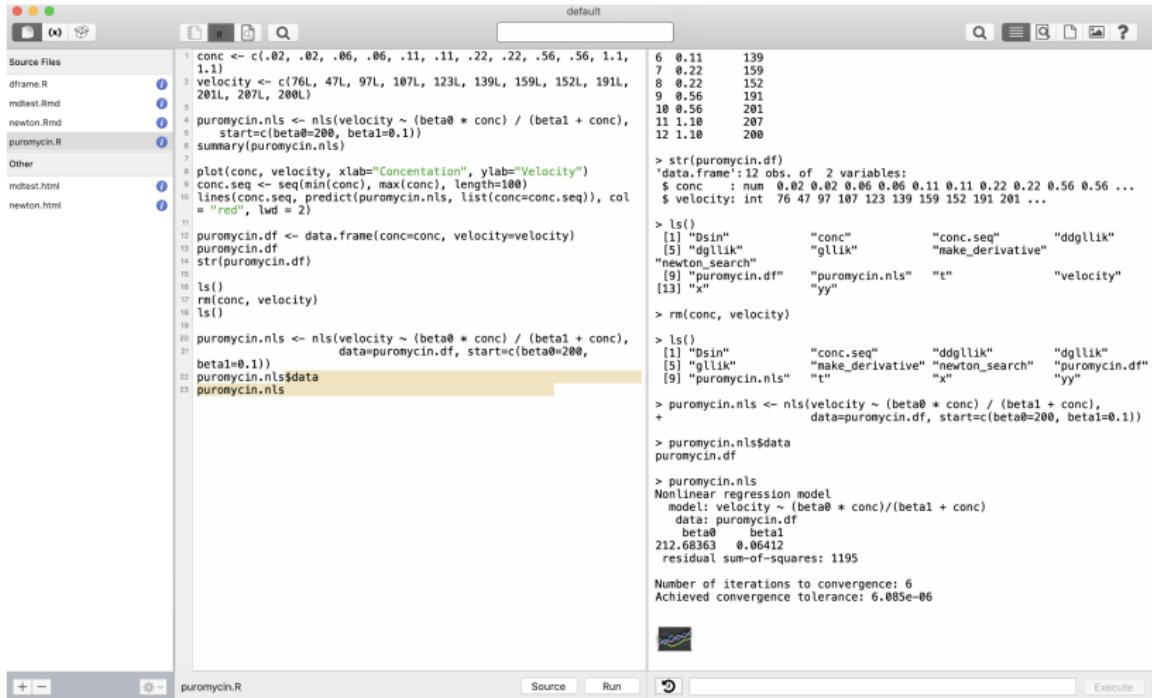
> puromycin.df <- data.frame(conc=conc, velocity=velocity)

> puromycin.df
  conc velocity
1  0.02      76
2  0.02      47
3  0.06      97
4  0.06     187
5  0.11     123
6  0.11     139
7  0.22     159
8  0.22     152
9  0.56     191
10 0.56     201
11 1.18     287
12 1.18     280
```

The results show the estimated parameters: beta0 = 2.127e+02 and beta1 = 6.412e-02. The plot shows a good fit of the model to the data points. The final data frame contains 12 rows of concentration and velocity values.

Figure 3: Rc² ClientR

Rc² Client with an R Image Icon



The screenshot shows the Rc² Client interface running on a Mac OS X system. The window title is "default". The left sidebar lists "Source Files" including dframe.R, rdistest.Rmd, newton.Rmd, and puromycin.R. The main pane displays R code for fitting a non-linear regression model to puromycin data. The code includes generating data frames, fitting models with nls(), and printing results. A small R logo icon is visible in the bottom right corner of the main pane.

```
1 conc <- c(.02, .02, .06, .06, .11, .11, .22, .22, .56, .56, 1.1, 1.1)
2 velocity <- c(76L, 47L, 97L, 107L, 123L, 139L, 159L, 152L, 191L, 261L, 267L, 260L)
3
4 puromycin.nls <- nls(velocity ~ (beta0 * conc) / (beta1 + conc),
5   start=c(beta0=200, beta1=0.1))
6 summary(puromycin.nls)
7
8 plot(conc, velocity, xlab="Concentration", ylab="Velocity")
9 conc.seq <- seq(min(conc), max(conc), length=100)
10 lines(conc.seq, predict(puromycin.nls, list(conc=conc.seq)), col = "red", lwd = 2)
11
12 puromycin.df <- data.frame(conc=conc, velocity=velocity)
13 puromycin.df
14 str(puromycin.df)
15
16 ls()
17 rm(conc, velocity)
18 ls()
19
20 puromycin.nls <- nls(velocity ~ (beta0 * conc) / (beta1 + conc),
21   data=puromycin.df, start=c(beta0=200,
22   beta1=0.1))
23 puromycin.nls$data
24 puromycin.nls
```

```
6 0.11    139
7 0.22    159
8 0.22    152
9 0.56    191
10 0.56    201
11 1.10    207
12 1.10    208

> str(puromycin.nls)
'data.frame': 12 obs. of  2 variables:
 $ conc : num  0.02 0.02 0.06 0.06 0.11 0.11 0.22 0.22 0.56 0.56 ...
 $ velocity: int  76 47 97 107 123 139 159 152 191 201 ...

> ls()
[1] "Dsin"           "conc"          "conc.seq"       "ddgllik"
[5] "dglilik"        "gllik"         "make_derivative"
[9] "newton_search"  "puromycin.df"  "t"              "velocity"

> rm(conc, velocity)

> ls()
[1] "Dsin"           "conc.seq"       "ddgllik"        "dglilik"
[5] "gllik"          "make_derivative" "newton_search"  "puromycin.df"
[9] "puromycin.nls"  "t"              "x"              "yy"

> puromycin.nls <- nls(velocity ~ (beta0 * conc) / (beta1 + conc),
+   data=puromycin.df, start=c(beta0=200, beta1=0.1))

> puromycin.nls$data
puromycin.nls

> puromycin.nls
Nonlinear regression model
  model: velocity ~ (beta0 * conc) / (beta1 + conc)
  data: puromycin.df
    beta0      beta1
  212.68363  0.06412
  residual sum-of-squares: 1195

Number of iterations to convergence: 6
Achieved convergence tolerance: 6.085e-06
```

Figure 4: Rc² ClientR2

Rc² Client with an R Plot

default

Source Files

- dframe.R
- ndtest.Rmd
- newton.Rmd
- puromycin.R**
- Other
- ndtest.html
- newton.html

```
1 conc <- c(.02, .02, .06, .06, .11, .11, .22, .22, .56, .56, 1.1,
1.1)
2 velocity <- c(76L, 47L, 97L, 107L, 123L, 139L, 159L, 152L, 191L,
281L, 267L, 280L)
3
4 puromycin.nls <- nls(velocity ~ (beta0 * conc) / (beta1 + conc),
5 start=c(beta0=200, beta1=0.1))
6 summary(puromycin.nls)
7
8 plot(conc, velocity, xlab="Concentration", ylab="Velocity")
9 conc.seq <- seq(min(conc), max(conc), length=100)
10 lines(conc.seq, predict(puromycin.nls, list(conc=conc.seq)), col =
"red", lwd = 2)
11
12 puromycin.df <- data.frame(conc=conc, velocity=velocity)
13 puromycin.df
14 str(puromycin.df)
15
16 ls()
17 rm(conc, velocity)
18 ls()
19
20 puromycin.nls <- nls(velocity ~ (beta0 * conc) / (beta1 + conc),
21 data=puromycin.df, start=c(beta0=200,
22 beta1=0.1))
23 puromycin.nls$data
24 puromycin.nls
```

Velocity

Concentration

Source Run < > img104.png

Figure 5: Rc² ClientR3

Rc² Client with an Rmd File

The screenshot shows the Rc² Client interface with the following details:

- Source Files:** A sidebar listing files: dframe.R, ndtest.Rmd, newton.Rmd (selected), puromycin.R, Other, ndtest.html, and newton.html.
- Code Editor:** The main area displays R code for Newton's method. The code defines a function `newton_search` that takes four arguments: `f`, `df`, `guess`, and `conv`. It uses a while loop to iteratively improve the guess until the absolute value of the function at the guess is less than `conv`. It also includes a local function `improve` to calculate the next guess based on the current guess and the derivative at that point.
- Text Block:** To the right of the code editor, there is explanatory text about the Newton search method and its implementation.
- Output Area:** Below the code editor, the output of the R code is shown in three sections:
 - newton_search:** The function definition.
 - newton_search(f = sin, df = cos, guess = 3):** The result is `#> [1] 3.142547`.
 - newton_search(f = sin, df = cos, guess = 3, conv = 0.000001):** The result is `#> [1] 3.141593`.
- Text Block:** Below the output, there is a note about the value of π .
- Bottom Navigation:** Buttons for Source and Run, and navigation arrows.
- File Status:** Shows 'newton.Rmd' is the active file.

Figure 6: Rc² ClientRmd

Rc² Client with a Help File

The screenshot shows the Rc² Client interface with a help file open for the `plot` function. The left pane displays the R code for a script named `puromycin.R`, which includes a call to `plot`. The right pane shows the help documentation for the `plot` function, which is described as a generic X-Y Plotting function. The documentation includes sections for Description, Usage, Arguments, and type, along with examples and notes.

```
default
plot (graphics)
Generic X-Y Plotting

Description
Generic function for plotting of R objects. For more details about the graphical parameter arguments, see par.
For simple scatter plots, plot.default will be used. However, there are plot methods for many R objects, including functions, data.frames, density objects, etc. Use methods\(plot\) and the documentation for these.

Usage
plot(x, y, ...)

Arguments
x
  the coordinates of points in the plot. Alternatively, a single plotting structure, function or any R object with a plot method can be provided.
y
  the y coordinates of points in the plot, optional if x is an appropriate structure.
...
  Arguments to be passed to methods, such as graphical parameters (see par). Many methods will accept the following arguments:
type
  what type of plot should be drawn. Possible types are
  • "p" for points,
  • "l" for lines,
  • "b" for both,
  • "c" for the lines part alone of "b",
  • "o" for both 'overplotted'.
```

Q Filter puromycin.R Source Run < > R: Generic X-Y Plotting

Figure 7: Rc² Help

Rc² Client with a Search File



The screenshot shows a search interface with a search bar containing "coordinates". Below the search bar, there's a section titled "Arguments" with the following details:

- x the **coordinates** of points in the plot. Alternatively, a single plotting structure, function or *any R object with a plot method* can be provided.
- y the y **coordinates** of points in the plot, *optional* if x is an appropriate structure.
- ... Arguments to be passed to methods, such as [graphical parameters](#) (see [par](#)). Many methods will accept the following arguments:

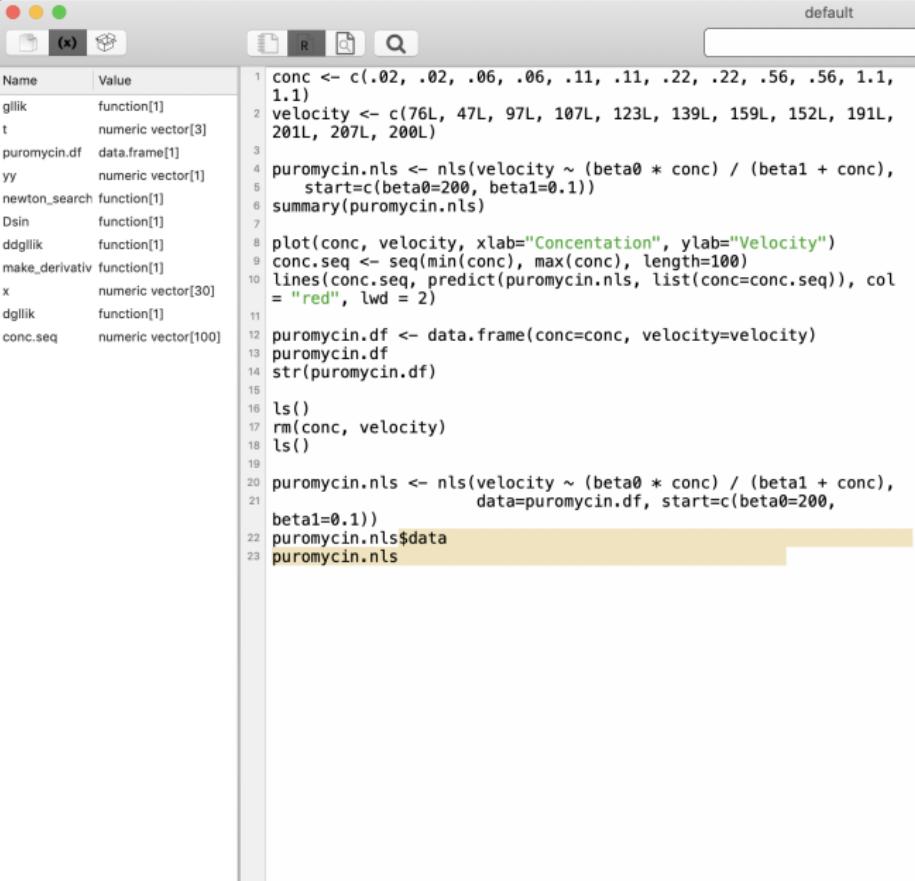
type

what type of plot should be drawn. Possible types are

- "p" for points,
- "l" for lines,
- "b" for both,
- "c" for the lines part alone of "b",
- "o" for both 'overplotted',
- "n" for 'histogram' like (or 'high-density') vertical lines,
- "s" for stair steps,
- "s" for other steps, see 'Details' below,
- "n" for no plotting.

All other types give a warning or an error; using, e.g., `type = "punkte"` being equivalent to `type = "p"` for S compatibility. Note that some methods, e.g. [plot.factor](#), do not accept this.

Rc² Client with an Environment Display



Name	Value
gllik	function[1]
t	numeric vector[3]
puromycin.df	data.frame[1]
yy	numeric vector[1]
newton_search	function[1]
Dsin	function[1]
ddgllik	function[1]
make_derivativ	function[1]
x	numeric vector[30]
dgllik	function[1]
conc.seq	numeric vector[100]

```
1 conc <- c(.02, .02, .06, .06, .11, .11, .22, .22, .56, .56, 1.1,
2 velocity <- c(76L, 47L, 97L, 107L, 123L, 139L, 159L, 152L, 191L,
3 201L, 207L, 200L)
4 puromycin.nls <- nls(velocity ~ (beta0 * conc) / (beta1 + conc),
5 start=c(beta0=200, beta1=0.1))
6 summary(puromycin.nls)
7
8 plot(conc, velocity, xlab="Concentration", ylab="Velocity")
9 conc.seq <- seq(min(conc), max(conc), length=100)
10 lines(conc.seq, predict(puromycin.nls, list(conc=conc.seq)), col
= "red", lwd = 2)
11
12 puromycin.df <- data.frame(conc=conc, velocity=velocity)
13 puromycin.df
14 str(puromycin.df)
15
16 ls()
17 rm(conc, velocity)
18 ls()
19
20 puromycin.nls <- nls(velocity ~ (beta0 * conc) / (beta1 + conc),
21 data=puromycin.df, start=c(beta0=200,
beta1=0.1))
22 puromycin.nls$data
23 puromycin.nls
```

Rc² Client with a Data Frame Display

The screenshot shows the Rc² Client interface with a data frame display. The interface has a top menu bar with icons for file operations and a title bar labeled "default". Below the title bar is a search bar. The main area contains three windows:

- puromycin.df**: A data frame window showing a table with two columns: "conc" and "velocity". The data consists of 12 rows of values.

	conc	velocity
1	0.02	76
2	0.02	47
3	0.06	97
4	0.06	107
5	0.11	123
6	0.11	139
7	0.22	159
8	0.22	152
9	0.56	191
10	0.56	201
11	1.1	207
12	1.1	200

- puromycin.nls**: A script editor window containing R code for non-linear least squares fitting. It includes a plot command with "conc" and "velocity" as arguments, and a call to "nls" with "start=c(beta0=200)".

```
16, .06, .11, .11, .22, .22, .56, .56, 1.1,  
.97L, 107L, 123L, 139L, 159L, 152L, 191L,  
  
velocity ~ (beta0 * conc) / (beta1 + conc),  
beta1=0.1))  
  
:lab="Concentration", ylab="Velocity")  
inc), max(conc), length=100)  
t(puromycin.nls, list(conc=conc.seq)), col  
  
frame(conc=conc, velocity=velocity)  
  
velocity ~ (beta0 * conc) / (beta1 + conc),  
data=puromycin.df, start=c(beta0=200,
```

Rc² Deployment

Kubernetes is available on AWS, Azure, and of course Google Cloud as a service. However, these cloud providers do not have predictable price models and users/developers can incur large unexpected costs.

Digital Ocean Kubernetes gives the user far more control over these costs, i.e., limits can be set and the user is only charged for what is actually used.