

How to run and track machine learning models in R with MLflow

This [article](#) has briefly described what [MLflow](#) is and how it works. *MLflow* currently provides APIs in Python language that users can invoke in their machine learning source codes to log parameters, metrics, and artifacts to be tracked by the *MLflow* tracking server.

Users familiar with R and perform machine learning operations in R may like to track their models and every runs with *MLflow*. There are several approaches users can take.

- Waiting for [MLflow](#) to release the APIs in R, or
- Wrapping *MLflow* RESTful APIs and logging through `curl` commands, or
- Calling existing Python APIs with some R packages that can invoke Python interpreter

The last approach is simple and easy enough while allows users to interact with *MLflow* without waiting for R APIs to be available. This tutorial will illustrate how to achieve this with R package [reticulate](#).

reticulate is an open source R package that allows to call Python from R by embedding a Python session within the R session. It provides seamless and high-performance interoperability between R and Python. The package is available in [CRAN repository](#).

MLflow also comes with a [Projects](#) component that packs data, source code with commands, parameters and execution environment setup together as a self-contained specification. Once a `MLproject` is defined, users can run it everywhere. Currently `MLproject` can run Python code or shell command. It can also set up the Python environment for the project specified in the `conda.yaml` file defined by users.

For R users, it is common to load some packages in the R source codes. These packages need to installed for the R code to run. In the future, it could be a good enhancement for *MLflow* to add something similar to `conda.yaml` to set up R package dependencies. This tutorial will show how to create a `MLproject` containing R source code and run it with `mlflow run` command.

Learning objectives

In this tutorial, developers will install and set up the *MLflow* environment, train and track machine learning models in R, package source codes and data in a `MLproject` and run with `mlflow run` command.

Prerequisites

Before beginning this tutorial, you should have Python installed on the platform where R is running. I prefer installing [miniconda](#). Since the machine learning training will be done in R, R should be already installed on the platform as well.

Estimated time

Completing this tutorial should take approximately 30 minutes.

Steps

Step 1: Install *MLflow*

Create a virtualenv for *MLflow* and install [mlflow](#) package as follow (with `conda`):

```
conda create -q -n mlflow python=3.6
source activate mlflow
pip install -U pip
pip install mlflow
```

Step 2: Install `reticulate` R package

Install [reticulate](#) package through R.

```
install.packages("reticulate")
```

R

`reticulate` allows R to call Python functions seamlessly. The Python package is loaded by the `import` statement. Calling to a function is through `$` operator.

```
> library(reticulate)
> path <- import("os.path")
> path$isdir("/tmp")
[1] TRUE
```

R

As you can see above, it is very simple to call Python functions in `os.path` module from R with this package. So you can do the same thing with `mlflow` package by importing it and then call

`mlflow$log_param` and `mlflow$log_metric` to log parameters and metrics for the R script.

Step 3: Train a GLM model with [SparkR](#)

Following R script builds a linear regression model with [SparkR](#). You need `SparkR` package installed for this [example](#).

```
# load the reticulate package and import mlflow Python module
library(reticulate)
mlflow <- import("mlflow")

# load SparkR package and start spark session
library(SparkR, lib.loc = c(file.path(Sys.getenv("SPARK_HOME"), "R", "lib")))
sparkR.session(master="local[*]")

# convert iris data.frame to SparkDataFrame
df <- as.DataFrame(iris)

# parameter for GLM
family <- c("gaussian")

# log the parameter
mlflow$log_param("family", family)

# fit the GLM model
model <- spark.glm(df, Species ~ ., family = family)

# exam the model
summary(model)

# path to save the model
model_path <- "/tmp/mlflow-GLM"

# save the model
write.ml(model, model_path)

# log the artifact
mlflow$log_artifacts(model_path)

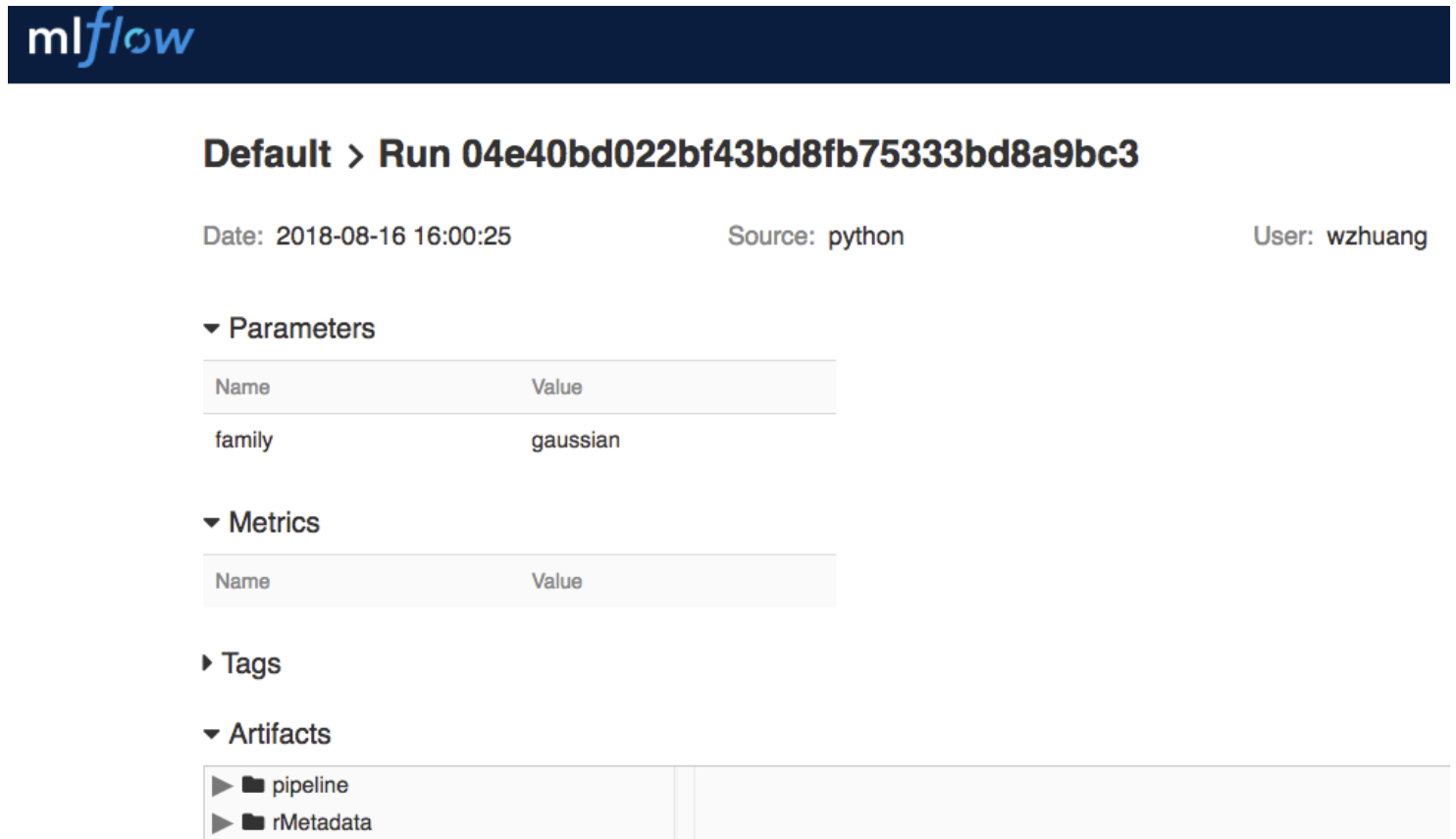
# stop spark session
sparkR.session.stop()
```

You can either copy the script to `R` or [Rstudio](#) and run interactively, or save it to a file and run with

`Rscript` command. Make sure that the `PATH` environment variable includes the path to the *mlflow* Python virtualenv.

Step 4: Launch the *MLflow* UI

Launch *MLflow* UI by running `mlflow ui` command from a shell. Then open browser and go to page link with url `http://127.0.0.1:5000`. Your previous GLM model training is now showing and so it can be tracked. Here is a snapshot.



The screenshot shows the MLflow UI interface. At the top is the MLflow logo. Below it, the breadcrumb 'Default > Run 04e40bd022bf43bd8fb75333bd8a9bc3' is displayed. The run details include: Date: 2018-08-16 16:00:25, Source: python, and User: wzhuang. There are four expandable sections: Parameters, Metrics, Tags, and Artifacts. The Parameters section is expanded, showing a table with one row: family (gaussian). The Artifacts section is also expanded, showing two items: pipeline and rMetadata, each with a folder icon and a play button.

Name	Value
family	gaussian

Name	Value
------	-------

Tags

Artifacts

- pipeline
- rMetadata

Step 5: Train a decision tree model

Download the [wine-quality.csv](#) data to be learned to your platform.

Install the `rpart` package on your R environment:

```
install.packages("rpart")
```

R

Follow this example [rpart-example.R](#) to fit a tree model:

```

# Source prep.R file to install the dependencies
source("prep.R")

# Import mlflow python package for tracking
library(reticulate)
mlflow <- import("mlflow")

# Load rpart to build a tree model
library(rpart)

# Read in data
wine <- read.csv("wine-quality.csv")

# Build the model
fit <- rpart(quality ~ ., wine)

# Save the model that can be loaded later
saveRDS(fit, "fit.rpart")

# Save the model to mlflow tracking server
mlflow$log_artifact("fit.rpart")

# Plot
jpeg("rplot.jpg")
par(xpd=TRUE)
plot(fit)
text(fit, use.n=TRUE)
dev.off()

# Save the plot to mlflow tracking server
mlflow$log_artifact("rplot.jpg")

```

The R code above includes three parts: the model training, the artifacts logging through *MLflow*, and the R package dependencies installation.

Step 6: Prepare package dependencies for MLproject

In above example, these two R packages, `reticulate` and `rpart`, are required for the code to run. To pack these codes into a self-contained project, some sort of script should be run to automatically install these packages if the platform does not have them installed.

Any specific R package needed for the project is going to be installed through `prep.R` with these codes:

```
# Accept parameters, args[6] is the R package repo url
args <- commandArgs()

# All installed packages
pkgs <- installed.packages()

# List of required packages for this project
reqs <- c("reticulate", "rpart")

# Try to install the dependencies if not installed
sapply(reqs, function(x){
  if (!x %in% rownames(pkgs)) {
    install.packages(x, repos=c(args[6]))
  }
})
```

Step 7: Test your codes

Before packaging these into a *MLproject*, try to test by directly invoking `Rscript` command as follow:

```
Rscript rpart-example.R https://cran.r-project.org/
```

From the *MLflow* UI, you should see this run been tracked like this screen snapshot:

Default > Run 5edf5d45812b48b3a4433995519a2553

Date: 2018-09-04 16:42:45

Source:  python

User: wzhuang

▼ Parameters

Name	Value
------	-------

▼ Metrics

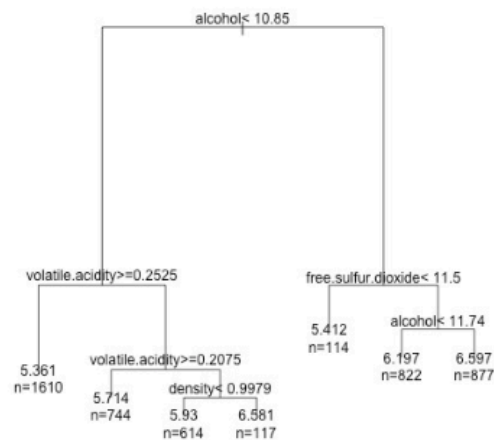
Name	Value
------	-------

► Tags

▼ Artifacts

-  fit.rpart
-  rplot.jpg

Full Path: /Users/wzhuang/DocsDump/files/mlflow/projects/R/mlruns/0/5edf5d45812b48b3a4433995519a2553/artifacts/rplot.jpg
Size: 22.9KB



Step 8: Create a MLproject

Now let's write the spec and pack this project into a *MLproject* that *MLflow* knows to run. All needed to be done is creating the [MLproject](#) file in the same directory.

```

name: r_example

entry_points:
  main:
    parameters:
      r-repo: {type: string, default: "https://cran.r-project.org/"}
      command: "Rscript rpart-example.R {r-repo}"
  
```

YAML

In this file, it defines a `r_example` project with a `main` entry point. The entry point specifies the command and parameters to be executed by the `mlflow run`. For this project, `Rscript` is the shell command to

invoke the R source code. `r-repo` parameter provides the URL string where the dependent packages can be installed from. A default value is set. This parameter is passed to the command running the R source code.

You have all the files required to train this tree model, you can create a `MLproject` by creating a directory and copying the data and R source codes to that directory.

```
.
├── R
│   ├── MLproject
│   ├── prep.R
│   ├── rpart-example.R
│   └── wine-quality.csv
```

Step 9: Check in and test the MLproject

The above `MLproject` can be checked in and pushed to github repository. To test the project, with following command, it can be run on any platform that has R installed.

```
mlflow run https://github.com/adrian555/DocsDump#files/mlflow-projects/R
```

The project can also be viewed from the *MLflow* tracking UI like this screen snapshot:

Default > Run c98e7e98972042469b3c1c750fe117d2

Date: 2018-09-05 10:21:21

Source:  R

Git Commit: b9607c

Entry Point: main

User: wzhuang

Duration: 6.1s

Run Command

```
mlflow run https://github.com/adrian555/DocsDump#files/mlflow-projects/R -v b9607cf73ae31b2369c94667cdc406059963af0e -P r-repo=https://cran.r-project.org/
```

▼ Parameters

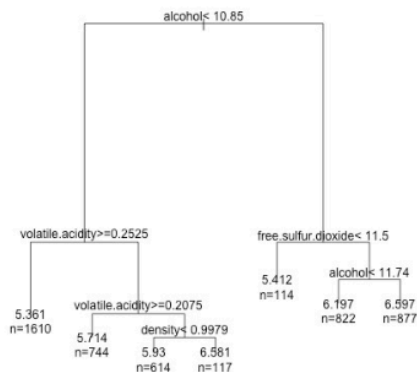
Name	Value
r-repo	https://cran.r-project.org/

▼ Metrics

Name	Value
------	-------

► Tags

▼ Artifacts

 fit.rpart
 rplot.jpgFull Path: /Users/wzhuang/mlruns/0/c98e7e98972042469b3c1c750fe117d2/artifacts/rplot.jpg
Size: 22.9KB

The differences between this view and previous run without `MLproject` spec are the `Run Command` which captures the exact command to run the project, and the `Parameters` which automatically logs any parameters passed to entry points.

Summary

In this tutorial, you have successfully created a `MLproject` in R, track and run it with *MLflow*. The approach taken here lets R users take benefit of *MLflow* `Tracking` component and track their R models in a quick way. It also demonstrates what `Projects` component of *MLflow* is designed for - to define the project and make it easily to be rerun. R users can quickly set up their projects and enjoy the easiness of tracking and running projects with *MLflow* once going through this tutorial.