

Machine Learning

Data Visualization

Date published: 2020-07-16

Date modified: 2022-04-11

CLOUDERA

<https://docs.cloudera.com/>

Legal Notice

© Cloudera Inc. 2023. All rights reserved.

The documentation is and contains Cloudera proprietary information protected by copyright and other intellectual property rights. No license under copyright or any other intellectual property right is granted herein.

Unless otherwise noted, scripts and sample code are licensed under the Apache License, Version 2.0.

Copyright information for Cloudera software may be found within the documentation accompanying each component in a particular release.

Cloudera software includes software from various open source or other third party projects, and may be released under the Apache Software License 2.0 (“ASLv2”), the Affero General Public License version 3 (AGPLv3), or other license terms. Other software included may be released under the terms of alternative open source licenses. Please review the license and notice files accompanying the software for additional licensing information.

Please visit the Cloudera software product page for more information on Cloudera software. For more information on Cloudera support services, please visit either the Support or Sales page. Feel free to contact us directly to discuss your specific needs.

Cloudera reserves the right to change any products at any time, and without notice. Cloudera assumes no responsibility nor liability arising from the use of products, except as expressly agreed to in writing by Cloudera.

Cloudera, Cloudera Altus, HUE, Impala, Cloudera Impala, and other Cloudera marks are registered or unregistered trademarks in the United States and other countries. All other trademarks are the property of their respective owners.

Disclaimer: EXCEPT AS EXPRESSLY PROVIDED IN A WRITTEN AGREEMENT WITH CLOUDERA, CLOUDERA DOES NOT MAKE NOR GIVE ANY REPRESENTATION, WARRANTY, NOR COVENANT OF ANY KIND, WHETHER EXPRESS OR IMPLIED, IN CONNECTION WITH CLOUDERA TECHNOLOGY OR RELATED SUPPORT PROVIDED IN CONNECTION THEREWITH. CLOUDERA DOES NOT WARRANT THAT CLOUDERA PRODUCTS NOR SOFTWARE WILL OPERATE UNINTERRUPTED NOR THAT IT WILL BE FREE FROM DEFECTS NOR ERRORS, THAT IT WILL PROTECT YOUR DATA FROM LOSS, CORRUPTION NOR UNAVAILABILITY, NOR THAT IT WILL MEET ALL OF CUSTOMER’S BUSINESS REQUIREMENTS. WITHOUT LIMITING THE FOREGOING, AND TO THE MAXIMUM EXTENT PERMITTED BY APPLICABLE LAW, CLOUDERA EXPRESSLY DISCLAIMS ANY AND ALL IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO IMPLIED WARRANTIES OF MERCHANTABILITY, QUALITY, NON-INFRINGEMENT, TITLE, AND FITNESS FOR A PARTICULAR PURPOSE AND ANY REPRESENTATION, WARRANTY, OR COVENANT BASED ON COURSE OF DEALING OR USAGE IN TRADE.

Contents

Built-in CML Visualizations.....	4
Simple Plots.....	4
Saved Images.....	4
HTML Visualizations.....	5
IFrame Visualizations.....	5
Grid Displays.....	6
Documenting Your Analysis.....	7
 Cloudera Data Visualization for ML.....	 7

Built-in CML Visualizations

You can use built-in CML tools to create data visualizations including simple plots, saved images, HTML and iFrame visualizations, and grid displays.

Simple Plots

Cloudera Machine Learning supports using simple plot to create data visualizations.

To create a simple plot, run a console in your favorite language and paste in the following code sample:

R

```
# A standard R plot
plot(rnorm(1000))
# A ggplot2 plot
library("ggplot2")
qplot(hp, mpg, data=mtcars, color=am,
facets=gear~cyl, size=I(3),
xlab="Horsepower", ylab="Miles per Gallon")
```

Python

```
import matplotlib.pyplot as plt
import random
plt.plot([random.normalvariate(0,1) for i in xrange(1,1000)])
```

Cloudera Machine Learning processes each line of code individually (unlike notebooks that process code per-cell). This means if your plot requires multiple commands, you will see incomplete plots in the workbench as each line is processed.

To get around this behavior, wrap all your plotting commands in one Python function. Cloudera Machine Learning will then process the function as a whole, and not as individual lines. You should then see your plots as expected.

Saved Images

You can display images within your reports.

Use the following commands:

R

```
library("cdsw")

download.file("https://upload.wikimedia.org/wikipedia/commons/2/29/Minard.
png", "/cdn/Minard.png")
image("Minard.png")
```

Python

```
import urllib
from IPython.display import Image
urllib.urlretrieve("http://upload.wikimedia.org/wikipedia/commons/2/29/Minar
d.png", "Minard.png")

Image(filename="Minard.png")
```

HTML Visualizations

Your code can generate and display HTML in Cloudera Machine Learning.

To create an HTML widget, paste in the following:

R

```
library("cdsw")
html('<svg><circle cx="50" cy="50" r="50" fill="red" /></svg>')
```

Python

```
from IPython.display import HTML
HTML('<svg><circle cx="50" cy="50" r="50" fill="red" /></svg>')
```

Scala

Cloudera Machine Learning allows you to build visualization libraries for Scala using [jvm-repr](#). The following example demonstrates how to register a custom HTML representation with the "text/html" mimetype in Cloudera Machine Learning. This output will render as HTML in your workbench session.

```
//HTML representation
case class HTML(html: String)
//Register a displayer to render html
Displayers.register(classOf[HTML],
  new Displayer[HTML] {
    override def display(html: HTML): java.util.Map[String, String] = {
      Map(
        "text/html" -> html.html
      ).asJava
    }
  })

val helloHTML = HTML("<h1> <em> Hello World </em> </h1>")

display(helloHTML)
```

Iframe Visualizations

Most visualizations require more than basic HTML. Embedding HTML directly in your console also risks conflicts between different parts of your code. The most flexible way to embed a web resource is using an [IFrame](#).



Note:

Cloudera Machine Learning versions 1.4.2 (and higher) added a new feature that allowed users to [HTTP security headers](#) for responses to Cloudera Machine Learning. This setting is enabled by default. However, the X-Frame-Options header added as part of this feature blocks rendering of iFrames injected by third-party data visualization libraries.

To work around this issue, a site administrator can go to the [Admin Security](#) page and disable the Enable HTTP security headers property. Restart Cloudera Machine Learning for this change to take effect.

R

```
library("cdsw")
iframe(src="https://www.youtube.com/embed/8pHzROPlD-w", width="854px", height="510px")
```

Python

```
from IPython.display import HTML
HTML('<iframe width="854" height="510" src="https://www.youtube.com/embed/8pHzROPlD-w"></iframe>')
```

You can generate HTML files within your console and display them in IFrames using the /cdn folder. The cdn folder persists and services static assets generated by your engine runs. For instance, you can embed a full HTML file with IFrames.

R

```
library("cdsw")
f <- file("/cdn/index.html")
html.content <- paste("<p>Here is a normal random variate:", rnorm(1), "</p>")
writeLines(c(html.content), f)
close(f)
iframe("index.html")
```

Python

```
from IPython.display import HTML
import random

html_content = "<p>Here is a normal random variate: %f </p>" % random.normalvariate(0,1)

file("/cdn/index.html", "w").write(html_content)
HTML("<iframe src=index.html>")
```

Cloudera Machine Learning uses this feature to support many rich plotting libraries such as htmlwidgets, Bokeh, and Plotly.

Grid Displays

Cloudera Machine Learning supports built-in grid displays of DataFrames across several languages.

Python

Using DataFrames with the pandas package requires per-session activation:

```
import pandas as pd
pd.DataFrame(data=[range(1,100)])
```

For PySpark DataFrames, use pandas and run `df.toPandas()` on a PySpark DataFrame. This will bring the DataFrame into local memory as a pandas DataFrame.

**Note:**

A Python project originally created with engine 1 will be running pandas version 0.19, and will not auto-upgrade to version 0.20 by simply selecting engine 2 in the project's Settings Engine page.

The pandas data grid setting only exists starting in version 0.20.1. To upgrade, manually install version 0.20.1 at the session prompt.

```
!pip install pandas==0.20.1
```

R

In R, DataFrames will display as grids by default. For example, to view the Iris data set, you would just use:

```
iris
```

Similar to PySpark, bringing Sparklyr data into local memory with `as.data.frame` will output a grid display.

```
sparkly_df %>% as.data.frame
```

Scala

Calling the `display()` function on an existing dataframe will trigger a collect, much like `df.show()`.

```
val df = sc.parallelize(1 to 100).toDF()  
display(df)
```

Documenting Your Analysis

Cloudera Machine Learning supports Markdown documentation of your code written in comments.

This allows you to generate reports directly from valid Python and R code that runs anywhere, even outside Cloudera Machine Learning. To add documentation to your analysis, create comments in [Markdown](#) format:

R

```
# Heading  
# -----  
#  
# This documentation is important.  
#  
# Inline math:  $e^x$   
#  
# Display math:  $y = \sigma x + \epsilon$   
  
print("Now the code!")
```

Python

```
# Heading  
# -----  
#  
# This documentation is important.  
#  
# Inline math:  $e^x$   
#  
# Display math:  $y = \sigma x + \epsilon$   
  
print("Now the code!")
```

Cloudera Data Visualization for ML

Cloudera Data Visualization enables you to explore data and communicate insights across the whole data lifecycle by using visual objects. The fast and easy self-service data visualization streamlines collaboration in data analytics through the common language of visuals.

Using this rich visualization layer enables you to accelerate advanced data analysis. The web-based, no-code, drag-and-drop user interface is highly intuitive and enables you to build customized visualizations on top of your datasets,

build dashboards and applications, and publish them anywhere across the data lifecycle. This solution allows for customization and collaboration, and it provides you with a dynamic and data-driven insight into your business.

In CDP Public Cloud, Data Visualization is integrated with Cloudera Machine Learning (CML). For on prem use, Data Visualization is integrated with Cloudera Data Science Workbench (CDSW) workflows. You can use the same visualization tool for structured, unstructured/text, and ML analytics, which means deeper insights and more advanced dashboard applications. You can create native data visualizations to provide easy predictive insights for business users and accelerate production ML workflows from raw data to business impact.

For more information, see the *Cloudera Data Visualization documentation*.

Related Information

[Cloudera Data Visualization in CDP Public Cloud](#)

[Cloudera Data Visualization in CDSW](#)