

Cloudera Data Science Workbench

Lab guide

CLOUDERA

Note: Only use this document during the sessions. Copy or reproduce the content of this document is not permitted.

Cloudera Data Science Workbench Labs

Table of Contents

Introduction	3
Lab 1 - Login to Cloudera Data Science Workbench (CDSW)	4
Lab 2 - Creating a new project	5
Lab 3 - Visualization and Sharing	14
Lab 4 - Hadoop Integration	18
Lab 5 – Working with Models	20
Lab 6 - SparklyR	29
Lab 7 - Shiny	32
Lab 8 - Scala	34
Lab 9 - Scheduling Jobs	38
Challenge Lab 10 – Experiments	46
Challenge Lab 11 - Pushing the Boundaries	53
Appendix	55
Cloudera Documentation	55
CDSW User Guide	55

Introduction

Cloudera Data Science workbench is a new product from Cloudera launched in May 2017. It is based on the acquisition of Sense.io that we made in March 2016. Cloudera has taken this product enhanced it and ensures that all workloads can be pushed down to Cloudera.

Accelerate data science from exploration to production using R, Python, Spark and more

For data scientists



Open data science, your way.

Use R, Python, or Scala with your favorite libraries and frameworks



No need to sample.

Directly access data in secure Hadoop clusters through Apache Spark and Apache Impala



Reproducible, collaborative research.

Share insights with your whole team

For IT professionals



Bring analysis to the data.

Give your data science team the freedom to work how they want, when they want



Secure by default.

Stay compliant with out-of-the-box support for full Hadoop security



Flexible deployment.

Run on-premises or in the cloud

Cloudera Data Science workbench supports the R, Python, Scala programming languages. That capability could certainly be useful to Cloudera; the software could enable companies to make the most of their data scientists, who can then be more efficient with their use of company time and infrastructure.



Programming language and software environment for statistical computing and graphics.

Best known in: **Academia and statistics community.**



High-level programming language for general-purpose programming.

Best known in: **Machine learning and data engineering community.**



General-purpose functional programming language with a strong static type system.

Best known in: **Data engineering community, due to Spark.**

Cloudera's goal with Cloudera Data Science workbench is to Help more data scientists use the power of Hadoop, make it easy and secure to add new users, use cases.

Why Hadoop for Data Science well here are the reasons:

High volume, low cost shared storage – More data more kinds of data
Parallel compute local to the data – more experiments, better results
Scalable, fault tolerant – easy to scale out, not just scale up
Flexible multipurpose data platform – easier path to production
Superior flexibility and price / performance to any other data platform

Lab 1 - Login to Cloudera Data Science Workbench (CDSW)

In this lab you'll learn how to:

- Login to a Cloudera Data Science Workbench instance
- Set your Hadoop Authentication
- Navigate the Cloudera Data Science Workbench application

First thing you need to do is register onto Cloudera Data Science Workbench

URL: <http://cdsw.scroll.cl/>

Please register yourself using the link "Sign Up for a New Account". After registration, tell to instructor the username you chose.

Lab 2 - Creating a new project

You will see the project window as follows

The screenshot shows the Cloudera Data Science Labs interface. At the top, there's a header with the user name 'amolenaar' and a 'Projects' tab. Below the header are four circular performance indicators: 'sessions running' (0), 'jobs running' (0), 'models running' (0), and two for 'vCPU' and 'B' (both at 0). A search bar says 'Project quick find' and a user icon says 'amolenaar'. On the left, a sidebar has icons for 'Projects' (selected), 'Sessions', 'Experiments', 'Models', 'Jobs', and 'Settings'. The main area shows a section for 'Projects' with a message: 'There are no projects created by amolenaar yet.' It also includes links to 'Create a Team' and 'Working with a team makes it easier to manage permissions and projects.' A blue 'New Project' button is located in the bottom right corner.

You have on the left hand panel

Projects - where you create data science projects

Jobs - Run and schedule jobs and add dependencies

Sessions - Python, Scala or R sessions

Experiments - batch experiments

Models - build, deploy, and manage models as REST APIs to serve predictions

Settings - User, Hadoop Authentication, SSH Keys and permission settings

In the top right hand corner you have

Search bar - for search for projects

+ adding new projects or new teams

User name - Account settings and Sign out - Same as settings in home screen

Let's create a new Project

New Project

Copy and paste this GitHub URL into the Git tab

https://github.com/andremolenaar/workshop_tellarius.git

Project_name = your user name and labs

You should see this

Create a New Project

Project Name

CDSW Training01

Project Visibility

- Private** - Only added collaborators can view the project.
- Public** - All authenticated users can view this project.

Initial Setup

Blank

Template

Local

Git

https://github.com/andremolenaar/workshop_tellarius.git

Create Project

The screenshot shows the Cloudera Data Science Labs interface. On the left is a vertical sidebar with icons for Account, Overview, Sessions, Experiments, Models, Jobs, Files, Team, and Settings. The main area shows a project titled "CDSW Demo Test". Under "Models", it says "This project has no models yet. Create a [new model](#)". Under "Jobs", it says "This project has no jobs yet. Create a [new job](#) to document your analytics pipelines". Under "Files", there is a list of files: data, 1_python.py, 2_pyspark.py, 3_tensorflow.py, 4_sparklyr.R, 5_shiny.R, README.md, server.R, spark-defaults.conf, ui.R, and utils.py.

On the left hand side panel you will see a new menu items, among them Team where you can add team members to your project. Ask your neighbor for his or her username, and start typing in the search box.

As an example

A screenshot of the "Team" section. A search bar at the top contains "fi". Below it is a table with three rows. The first row shows "Filippo (flambiente)" with an orange profile icon, and a blue "Add" button. The second row shows "Sofie (Gundersen)" with an orange profile icon. The third row shows "Sofia Thorén (sofiathoren)" with an orange profile icon. To the right of the table, under "Permission", it says "Admin".

You can add anybody to your project. If you cannot find your neighbor, you can use the 'admin' user to share your project with.

Cloudera - Data Science Labs

Collaborators

This project is **private**. Only collaborators can view and edit this project. [Change Settings](#).

Add Collaborator

Collaborator	Permission	
A amolenaar	Admin	change delete
A admin	Viewer	change delete

Granting write or admin permission to other users may have security impact since it gives them full access to your project files and running sessions. Write or admin permission should only be granted to trusted users.

Click on the ‘Overview’ button.

Now, Launch a Python 3 Session as shown below:

The screenshot shows the Cloudera Data Science Workbench interface for a project named "Andre_CDSW_Demos". The left sidebar includes links for Account, Overview, Sessions, Experiments, Models, Jobs, Files, Team, and Settings. A license notice at the bottom left states "License Expires in 17 days". The main content area displays the project overview, showing sections for Models (no models), Jobs (no jobs), and Files. The Files section lists several Python and R scripts, including "1_python.py", "2_pyspark.py", "3_tensorflow.py", "4_sparklyr.R", "5_shiny.R", "README.md", "server.R", "spark-defaults.conf", "ui.R", and "utils.py". The "Open Workbench" button in the top right corner is circled in red. Below the files list, there is a section titled "Cloudera Data Science Workbench demos" with a brief description and a link to "Basic Python visualizations (Python 2)".

Cloudera - Data Science Labs

The screenshot shows the Cloudera Data Science Workbench interface. On the left is a file browser with a dark theme, showing files like 1_python.py, 2_pyspark.py, 3_tensorflow.py, 4_sparklyr.R, 5_shiny.R, README.md, server.R, spark-defaults.conf, ui.R, and utils.py. In the center is a code editor window titled 'README.md' containing Python code. On the right is a 'Start New Session' panel. It includes sections for 'Engine Image - Configure' (Base Image v6 - docker.repository.cloudera.com/cdsw/engine:6), 'Select Engine Kernel' (Python 3 selected), 'Select Engine Profile' (1 vCPU / 4 GiB Memory, 0 GPUs), and two buttons: 'Launch Session' (circled in red) and 'Run Experiment..'. A refresh icon is also present.

1. On the far left is a file browser (note the little ‘refresh’ icon at the top: )
2. In the middle is an editor open on the file that you selected - in this case the README.md file.
3. On the right is a Session Start tile - in this case it’s waiting for you to select an engine to run (so far your project has file space but no compute sessions).
4. Select the Python 3 kernel; select the 1 vCPU/4GiB Engine (and use this size engine for all your Python sessions in this workshop) and then the Launch Session button.

Start New Session

Engine Image - Configure

Base Image v6 - docker.repository.cloudera.com/cdsw/engine:6

Select Engine Kernel

- Python 2
- Python 3
- Scala
- R

Select Engine Profile

1 vCPU / 4 GiB Memory  0 GPUs 

 Launch Session

or Run Experiment..

5. This will startup a Python engine and the right hand side will become two tiles. The top one is an output tile and the bottom one (with a red, then green left hand border) is a shell input window.

Cloudera - Data Science Labs

The screenshot shows a Jupyter Notebook interface. On the left, a sidebar lists files: CDSW Demo Test, 1.python.py, 2.pyspark.py, 3.tensorflow.py, 4.sparklyr.R, 5.shiny.R, data, README.md, server.R, spark-defaults.conf, ui.R, and utils.py. The README.md file is open, displaying its content. On the right, an 'Untitled Session' tab is active, showing a Python 2 session. The session output includes instructions for setting up Python and R environments, upgrading packages, and running a job. The status bar at the bottom indicates 53 Lines, Markdown, and Spaces 2.

```
1 # Cloudera Data Science Workbench demos
2 Basic tour of Cloudera Data Science Workbench.
3
4 ## Workbench
5 There are 4 scripts provided which walk through the interactive capabilities
6
7 1. **Basic Python visualizations (Python 2).** Demonstrates:
8   - Markdown via comments
9   - Jupyter-compatible visualizations
10  - Simple console sharing
11 2. **PySpark (Python 2).** Demonstrates:
12   - Access to Hadoop HDFS CLI (e.g. `hdfs dfs -ls /`).
13 3. **Tensorflow (Python 2).** Demonstrates:
14   - Ability to install and use custom packages (e.g. `pip search tensorflow`)
15 4. **R on Spark via Sparklyr (R).** Demonstrates:
16   - Use familiar dplyr with Spark using [Sparklyr](http://spark.rstudio.com)
17 5. **Advanced visualization with Shiny (R).** Demonstrates:
18   - Use of 'shiny' to provide interactive graphics inside CDSW
19
20
21 ## Jobs
22 We recommend setting up a **Nightly Analysis** job to illustrate how data
23
24
25 ## Setup instructions
26 Note: You only need to do this once.
27
28 1. In a Python 2 Session:
29   Python
30   !pip2 install --upgrade dask
31   !pip2 install --upgrade keras
32   !pip2 install --upgrade matplotlib==2.0.0.
33   !pip2 install --upgrade pandas_highcharts
34   !pip2 install --upgrade protobuf
35   !pip2 install --upgrade tensorflow==1.3.0.
36   !pip2 install --upgrade seaborn
37
38 Note, you must then stop the session and start a new Python session in order
39
40 2. In an R Session:
41
42   install.packages('sparklyr')
43   install.packages('plotly')
44   install.packages("nycflights13")
45   install.packages("Lahman")
46   install.packages("mgcv")
47   install.packages('shiny')
48 ...
49
50 3. Stop all sessions, then proceed.
51
52
53
```

6. Look at the line 30 to 36 of README.md:

```
!pip2 install --upgrade dask
!pip2 install --upgrade keras
!pip2 install --upgrade matplotlib==2.0.0.
!pip2 install --upgrade pandas_highcharts
!pip2 install --upgrade protobuf
!pip2 install --upgrade tensorflow==1.3.0.
!pip2 install --upgrade seaborn
```

7. Question: What does this do? If you know Python it should be obvious; if you don't, then let me tell you: this is an execution of pip (a Python package manager), which will tell the system to install or upgrade the listed python packages
8. Make a block selection of these lines, and select 'Run Line(s)':

```

28 1. In a Python 2 Session:
29 ````Python
30 !pip2 install --upgrade dask
31 !pip2 install --upgrade keras
32 !pip2 install --upgrade matplotlib==2.0.0.
33 !pip2 install --upgrade pandas_highcharts
34 !pip2 install --upgrade protobuf
35 !pip2 install --upgrade tensorflow==1.3.0.
36 !pip2 install --upgrade seaborn
37
38 Note, you must then stop the session and start
39
40 2. In an R Session:
41 ````R

```

Run Line(s)

⌘Enter

Select All

⌘A

9. The cursor on the left should turn to red and, after a few seconds, you should see output like this:

The screenshot shows a Jupyter Notebook interface with two panes. The left pane contains a file named 'README.rst' with Python code for setting up a Data Science Workbench. The right pane shows the terminal output of running the code, which includes pip upgrade commands and dependency resolution logs.

```

1 # Cloudera Data Science Workbench demos
2 Basic tour of Cloudera Data Science Workbench.
3
4 ## Workbench
5 There are 4 scripts provided which walk through the interactive capabilities
6
7 1. **Basic Python visualizations (Python 2).** Demonstrates:
8   - Markdown via comments
9   - Jupyter-compatible visualizations
10  - Simple console sharing
11 2. **PySpark (Python 2).** Demonstrates:
12   - Execution of PySpark (kerberized) in YARN client mode.
13   - Access to Hadoop HDFS CLI (e.g. `hdfs dfs -ls /`).
14 3. **Tensorflow (Python 2).** Demonstrates:
15   - Ability to install and use custom packages (e.g. `pip search tensorflow`)
16 4. **R on Spark via Sparklyr (R).** Demonstrates:
17   - Use familiar dplyr with Spark using [Sparklyr](http://spark.rstudio.com)
18 5. **Advanced visualization with Shiny (R).** Demonstrates:
19   - Use of 'shiny' to provide interactive graphics inside CDSW
20
21 ## Jobs
22 We recommend setting up a **Nightly Analysis** job to illustrate how data
23
24
25 ## Setup instructions
26 Note: You only need to do this once.
27
28 1.. In a Python 2 Session:
29 ````Python
30 !pip2 install --upgrade dask
31 !pip2 install --upgrade keras
32 !pip2 install --upgrade matplotlib==2.0.0.
33 !pip2 install --upgrade pandas_highcharts
34 !pip2 install --upgrade protobuf
35 !pip2 install --upgrade tensorflow==1.3.0.
36 !pip2 install --upgrade seaborn
37
38 Note, you must then stop the session and start a new Python session in order
39
40 2. In an R Session:
41 ````R
42 install.packages('sparklyr')
43 install.packages('dplyr')
44 install.packages("nycflights13")
45 install.packages("Lahman")
46 install.packages("mgcv")
47 install.packages('shiny')
48 ...
49
50 3. Stop all sessions, then proceed.
51
52
53

```

Terminal Output:

```

.0.0 por-4.2.0 TENSORLOW-1.5.0 TENSORLOW-TENSORBOARD-0.1.8
You are using pip version 10.0.1, however version 18.0 is available.
You should consider upgrading via the 'pip install --upgrade pip' command.

> !pip2 install --upgrade seaborn
Collecting seaborn
  7251 Downloading https://files.pythonhosted.org/packages/7a/bf/04cfcc9616cedd4b5dd24dfc4039596ea9f50c1
db0d3f3e52b050f74a5/seaborn-0.9.0.tar.gz (198kB)
K 100% [=====] 204kB 4.3MB/s
?25hRequirement not upgraded as not directly required: numpy>=1.9.3 in /usr/local/lib/python2.7/site-packages (from seaborn) (1.12.1)
Requirement not upgraded as not directly required: scipy>=0.14.0 in /usr/local/lib/python2.7/site-packages (from seaborn) (1.1.0)
Requirement not upgraded as not directly required: pandas>=0.15.2 in /usr/local/lib/python2.7/site-packages (from seaborn) (0.20.1)
Requirement not upgraded as not directly required: matplotlib>=1.4.3 in /usr/local/lib/python2.7/site-packages (from seaborn) (2.0.0)
Requirement not upgraded as not directly required: python-dateutil in /usr/local/lib/python2.7/site-packages (from pandas>=0.15.2->seaborn) (2.7.3)
Requirement not upgraded as not directly required: pytz>=2011k in /usr/local/lib/python2.7/site-packages (from pandas>=0.15.2->seaborn) (2018.4)
Requirement not upgraded as not directly required: subprocess32 in /usr/local/lib/python2.7/site-packages (from matplotlib>=1.4.3->seaborn) (2.2.0)
Requirement not upgraded as not directly required: six>=1.10 in /usr/local/lib/python2.7/site-packages (from matplotlib>=1.4.3->seaborn) (1.11.0)
Requirement not upgraded as not directly required: functools32 in /usr/local/lib/python2.7/site-packages (from matplotlib>=1.4.3->seaborn) (3.2.3.post2)
Requirement not upgraded as not directly required: cyclere>=0.10 in /usr/local/lib/python2.7/site-packages (from matplotlib>=1.4.3->seaborn) (0.10.0)
Building wheels for collected packages: seaborn
  Running setup.py bdist_wheel for seaborn ... ?25ldone
?25h Stored in directory: /home/cdsweb/.cache/pip/wheels/fc/1c/74/c8f80a532c06a789599b8659b117ec7d7574cac4
a66f7dabfe
Successfully built seaborn
grin 1.2.1 requires argparse>=1.1, which is not installed.
bokeh 0.12.19 has requirement futures>=3.0.3, but you'll have futures 2.1.4 which is incompatible.
Installing collected packages: seaborn
Successfully installed seaborn-0.9.0
You are using pip version 10.0.1, however version 18.0 is available.
You should consider upgrading via the 'pip install --upgrade pip' command.

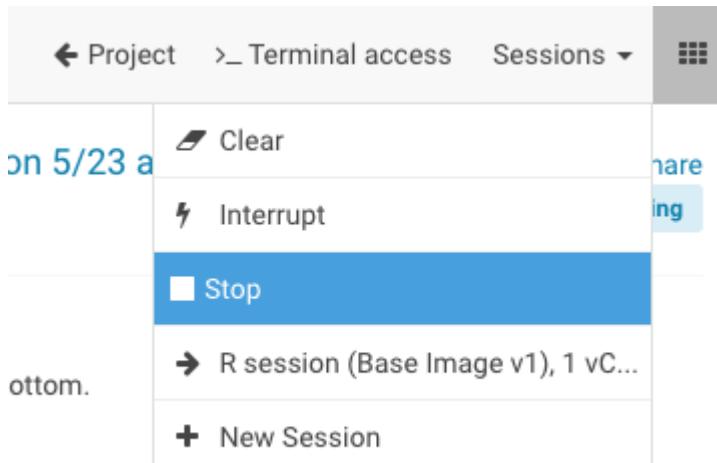
```

10. What's happening here is that the code in the file is being executed in the console (did you notice the left hand edge turned red?) and now you can see the output in the right hand screen.

11. Stop this session. The Stop button is either on the top menu bar (when there's sufficient room for it):

← Project >_ Terminal access ⚡ Clear ⚡ Interrupt █ Stop Sessions ▾

12. Or it's in the the Session drop down (when there's little room for buttons):



ter on Mac or **Ctrl-Enter** on Windows. You can also

13. Now launch an R session. **Use a 1vCPU, 4GiB Engine**. Use this size engine for **all** your R sessions during this workshop.

Start New Session

Engine Image - Configure

Base Image v6 - docker.repository.cloudera.com/cdsweb/engine:6

Select Engine Kernel

- Python 2
- Python 3
- Scala
- R

Select Engine Profile

1 vCPU / 4 GiB Memory ▾

0 GPUs ▾

Launch Session

or

Run Experiment..

14. This time we're going to execute lines 42 through 47 in the R session (these numbers could be off by 1 or so ... look at the images below and figure out what you need to select). Select and highlight just these lines then select Run Lines to execute as in prior step:

```
40 2. In an R Session:  
41   ```R  
42 install.packages('sparklyr')  
43 install.packages('plotly')  
44 install.packages("nycflights13")  
45 install.packages("Lahman")  
46 install.packages("mgcv")  
47 install.packages('shiny')
```

```
install.packages('sparklyr')  
install.packages('plotly')  
install.packages("nycflights13")  
install.packages("Lahman")  
install.packages("mgcv")  
install.packages('shiny')  
***
```

Run Line(s)

⌘ Enter

Select All

⌘ A

This process will take up to 10 minutes because code is being downloaded, compiled and installed into your project's workbench.

Note that this workbench is independent of any other project's workbench. You want to try out different and conflicting libraries? Go for it. Just start another project, open the workbench, pick the libraries you want, install them and off you go. Just like on your laptop, but this is managed, secure, stable, won't get stolen from a taxi, is always available and easily shared!

All of this isolation is achieved by mounting filesystems (one or more per project) into docker containers (one per engine). The details don't matter, except to note that now you can really go mad and try out lots of different and conflicting permutations without having to go down the complex path of virtual environments etc. It's all been done for you!

Lab 3 - Visualization and Sharing

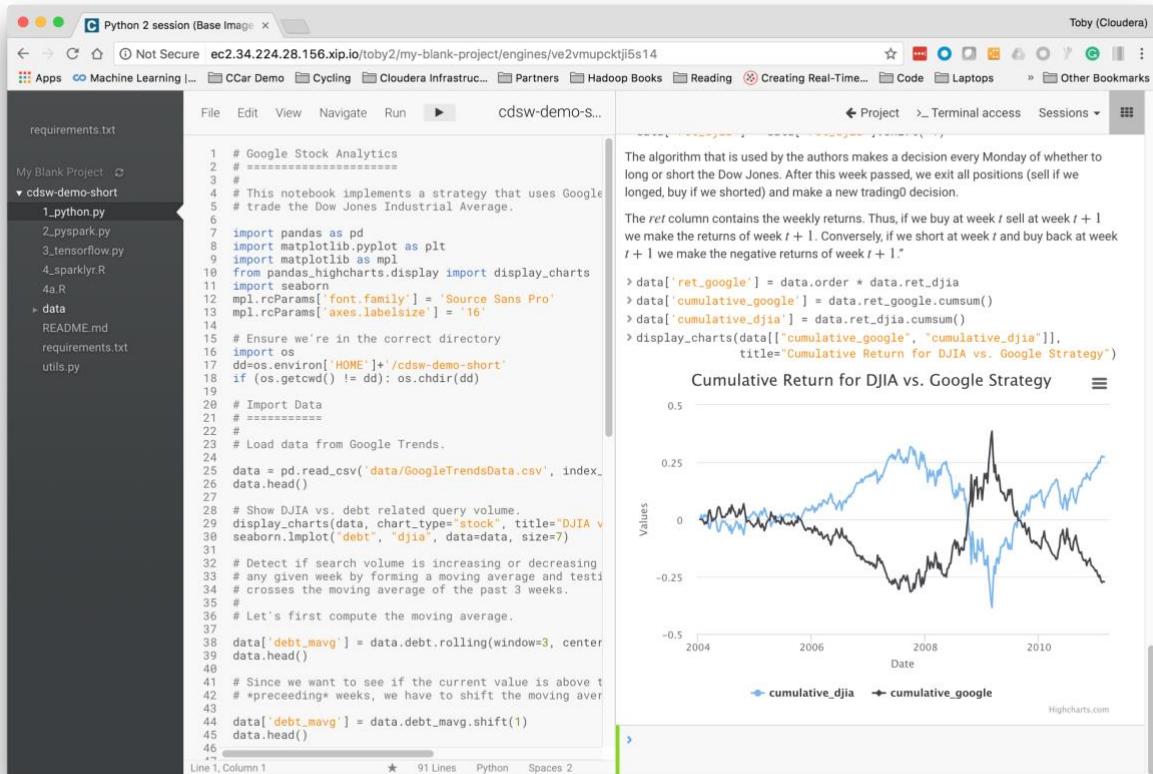
Check that you have NO sessions running - if any sessions are running then stop them now. You've setup your environment and those sessions can be safely disposed of.

Data Science is often about visualizing ideas, and then sharing them to persuade others to take action. CDSW lets you use the visualization tools you'd use naturally, and adds a neat twist to the whole idea of sharing. Let's get started:

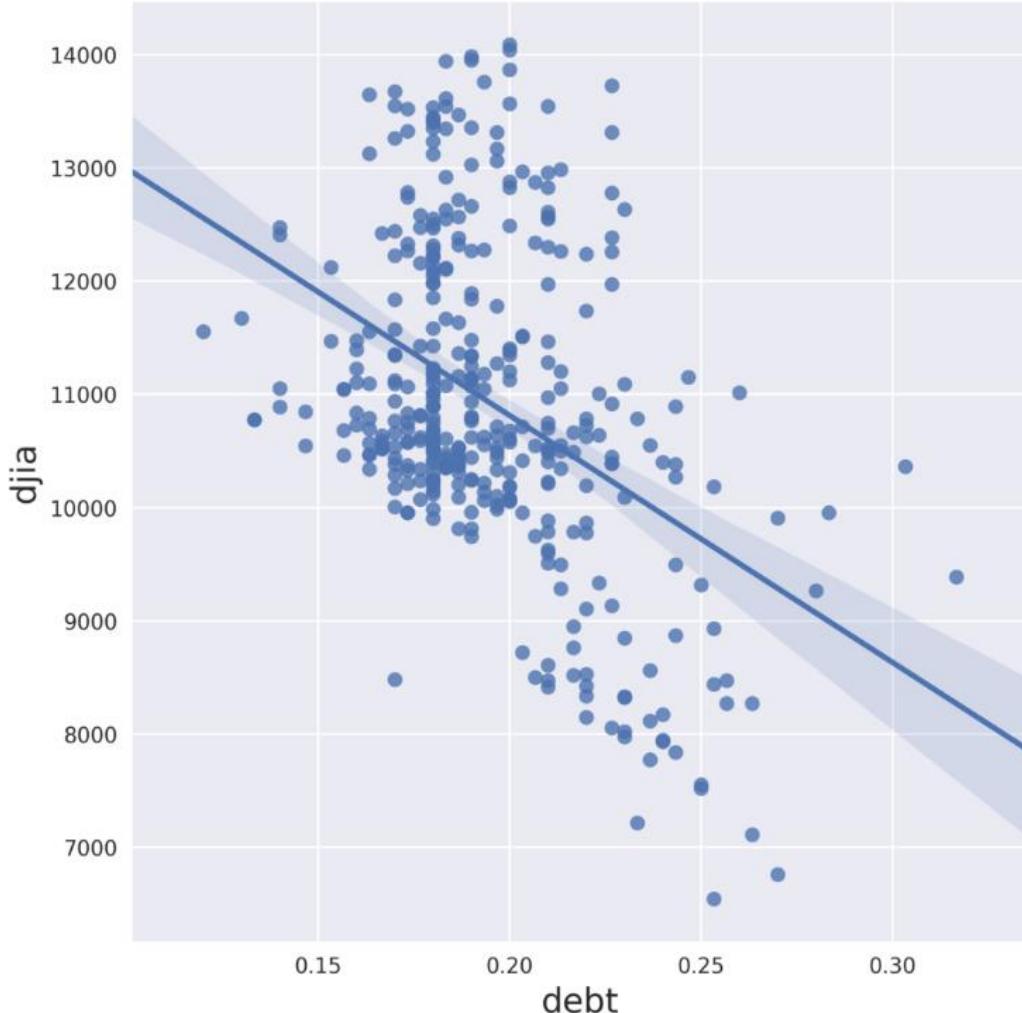
1. Start up a new Python 3 session (1vCPU, 4GiB) in the same manner you did before.
2. Select `1_python.py` in the file browser
3. Run the entire file (multiple ways of doing that - try to figure out more than one way. It should be pretty obvious!).



4. You should end up with some nice graphs in the output window:



5. You can see that CDSW is very similar to a notebook, supporting the same visualization tools. However, unlike a notebook, it doesn't use cells: instead it uses markup in the source file, and an output window. Furthermore, that window has some interesting properties ...
6. Scroll up to find this diagram:



7. On the left is a little chain link button:
8. Click on it and you'll see beneath the chart some html that can be used to embed that chart into a web site:

Copy and paste to embed this cell in your website!

width (optional)

height (optional)

```
<div id='sense-widget-s7dkodjq' class="sense-embed-container" style="width:100%;"><script id='embedding-script-21' src='http://consoles.cds.52.214.150.212.nip.io/js/sense-embed.js' data-cell-url='http://livelog.cds.52.214.150.212.nip.io/topics/xeupfj2oeyaktdpa_output/21' data-auth-token='eyJ0eXAiOiJKV1QiLCJhbGciOiJIUzI1NiJ9.eyJhdXRoijp7InJlYWQiOlsieGV1cGZqMm9leWFrdGRwYV8qIl0sIndyaXRlIjpbdata-assets-cdn-root='http://consoles.cds.52.214.150.212.nip.io/0/7/xeupfj2oeyaktdpa/' data-widget='sense-widget-s7dkodjq' data-width="" data-height=""></script></div>
```

9. Scroll to the top of the window and you'll see this on the far right (the exact layout depends upon the real estate available – you might have to expand your browser window to see the following links and they might be laid out vertically or horizontally):



10. Hit the 'collapse' link and see the difference in the output window.
11. Question: What difference did you see? How might this be used? Is it useful?
12. Notebooks have great output, but how do you share what they show you? CDSW solves this by simply providing a link to the output that you can send to anyone and they can see the output. Try it:

13. Select the 'Share' link:

These results are **private**. Only project members can view.

Share with Others

14. And then 'Share with Others' (your URL will be similar, but different from this one):

These results are being **shared**.

<http://cdsw.34.254.254.246.r>

Stop Sharing

Hide code and text

Who can view:

- Anonymous visitors with the link
- Any logged in user with the link
- Specific users/teams with the link ([Change...](#))
Currently shared with no one.

15. Cut and paste that link and put it into some other browser (best to be a completely different browser or an incognito windows than the one you're logged in with, but not that important)
16. You should see that you have access to almost the same output window (this new one doesn't have this share link!)

Python 2 session (Base Image v2), 1 vCPU / 2 GiB Memory, on 9/6 at 9:26

By Colm — Python 2 Session (Base Image v2) — 4 minutes ago for running

Google Stock Analytics

This notebook implements a strategy that uses Google Trends data to trade the Dow Jones Industrial Average.

Ensure we're in the correct directory

Import Data

Load data from Google Trends.

	djia	debt
Date		
2004-01-14	10485.18	0.210000
2004-01-22	10528.66	0.210000
2004-01-28	10702.51	0.210000
2004-02-04	10499.18	0.213333
2004-02-11	10579.03	0.200000

Show DJIA vs. debt related query volume.



So we've demonstrated how CDSW is like a notebook, but is perhaps more powerful, and has great sharing capability. Let's go on to see about integration with Hadoop!

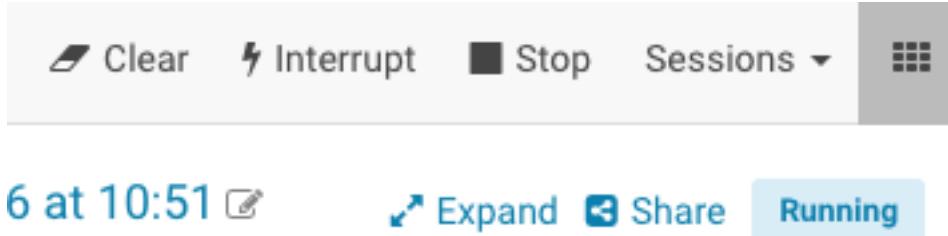
Lab 4 - Hadoop Integration

In this lesson we'll see two mechanisms for integrating with Hadoop:

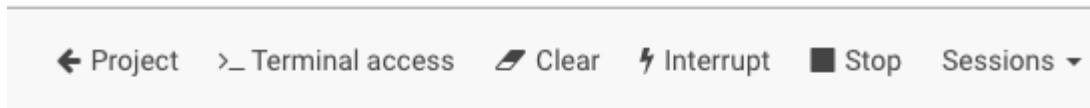
1. Filesystem - storing data in Hadoop itself using HDFS
2. Computation - executing code on the Hadoop cluster via Spark

Execute the following instructions.

1. Clean up your Python session by hitting the 'clear' button and the 'expand' link (if available)



2. Select the 2_pyspark.py file
3. If you don't have a Python 3 workbench session open yet, Launch a Python 3 session 1 vCPU / 4 GB RAM
4. Run line 34:
5. `!hdfs dfs -put -f $HOME/data/kmeans_data.txt /user/$HADOOP_USER_NAME`
6. Execute the 2_pyspark.py file in your already running Python session
7. **Question:** What did it do?
8. **Question:** What kind of thing is the variable 'data'? (try typing 'data' into the console and seeing what gets printed out.)
9. Open a terminal using the 'terminal' icon in the top right:



10. Execute 'hdfs dfs -ls' to see the data file in the hadoop file system (or, to show off, execute '! hdfs dfs -ls' in the python console to do the same thing!)

```
Welcome to Cloudera Data Science Workbench

Kernel: python3

Project workspace: /home/cdsd

Kerberos principal: training10@CLOUDERA.TELLARIUS.EU

Runtimes:
R: R version 3.5.1 (--) -- "Feather Spray"
Python 2: Python 2.7.11
Python 3: Python 3.6.1
Java: java version "1.8.0_131"

Git origin: https://github.com/andremolenaar/workshop_tellarius.git

cdsw@ck44c57bnm5xvzq2:~$ hdfs dfs -ls
Found 2 items
drwxrwx---  - training10 training10          0 2019-04-10 15:33 .sparkStaging
-rw-rw----  3 training10 training10        71 2019-04-10 15:33 kmeans_data.txt
cdsw@ck44c57bnm5xvzq2:~$ █
```

or

```
» !hdfs dfs -ls
Found 2 items
drwxrwx---  - training10 training10          0 2019-04-10 15:33 .sparkStaging
-rw-rw----  3 training10 training10        71 2019-04-10 15:33 kmeans_data.txt
```

If everything went correctly you'll see that we demonstrate:

- Natural integration with HDFS - it's just a path to a file!
- Natural parallel computation across the cluster using Spark

Lab 5 – Working with Models

Starting with version 1.4, Cloudera Data Science Workbench allows data scientists to build, deploy, and manage models as REST APIs to serve predictions.

Challenge

Data scientists often develop models using a variety of Python/R open source packages. The challenge lies in actually exposing those models to stakeholders who can test the model. In most organizations, the model deployment process will require assistance from a separate DevOps team who likely have their own policies about deploying new code.

For example, a model that has been developed in Python by data scientists might be rebuilt in another language by the devops team before it is actually deployed. This process can be slow and error-prone. It can take months to deploy new models, if at all. This also introduces compliance risks when you take into account the fact that the new re-developed model might not be even be an accurate reproduction of the original model.

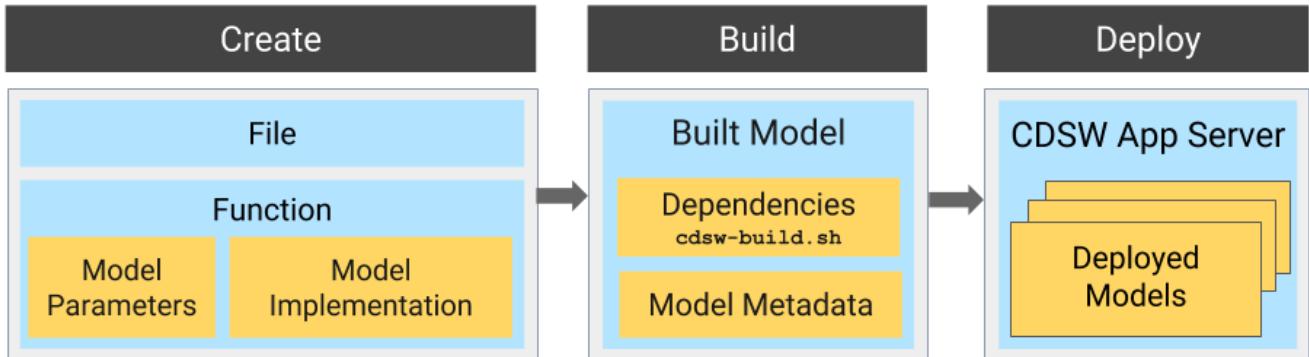
Once a model has been deployed, you then need to ensure that the devops team has a way to rollback the model to a previous version if needed. This means the data science team also needs a reliable way to retain history of the models they build and ensure that they can rebuild a specific version if needed. At any time, data scientists (or any other stakeholders) must have a way to accurately identify which version of a model is/was deployed.

Solution

Starting with version 1.4, Cloudera Data Science Workbench allows data scientists to build and deploy their own models as REST APIs. Data scientists can now select a Python or R function within a project file, and Cloudera Data Science Workbench will:

- Create a snapshot of model code, model parameters, and dependencies.
- Package a trained model into an immutable artifact and provide basic serving code.
- Add a REST endpoint that automatically accepts input parameters matching the function, and that returns a data structure that matches the function's return type.
- Save the model along with some metadata.
- Deploy a specified number of model API replicas, automatically load balanced.

Stages of the Model Deployment Process



Step 1: Examine the program `predic_churn_sklearn.py`

Open the project you created in the previous lab, and examine the file.

```
1 import pickle
2 import numpy as np
3
4 model = pickle.load(open("models/sklearn_rf.pkl", "rb"))
5
6 def predict(args):
7     account=np.array(args["feature"].split(", ")).reshape(1, -1)
8     return {"result" : model.predict(account)[0]}
9
10
```

This PySpark program uses the `pickle.load` mechanism to deploy models.. The model it refers to the `sklearn_rf.pkl` file, was saved in the previous lab from the experiment with the best predictive model.

There is a `predict` definition which is the function that calls the model, using features, and will return a `result` variable.

Step 2: Deploy the model

From the projects page of your project, select the ‘Models’ button.

The screenshot shows the 'Experiments and Models' page for a user named 'andre'. On the left is a sidebar with icons for Overview, Sessions, Experiments, Models (which is highlighted with a red circle), and Jobs. The main area displays the 'Models' section, which is currently empty. It includes a message: 'This project has no models yet. Create a [new model](#)'. Below this is the 'Jobs' section, also empty, with a message: 'This project has no jobs yet. Create a [new job](#) to document your analytics pipelines.' There is a '...' button at the bottom.

Select ‘New Model’, and populate specify the following configuration:

Name: something like “My Churn Prediction Model”
Description: Anything you want
File: predict_churn_sklearn.py
Function: predict
Example Inp: {
 "feature": "0, 65, 0, 137, 21.95, 83, 19.42, 111, 9.4, 6, 3.43, 4"
 }
Kernal: Python 3
Engine: 0.5 vCPU / 2 GiB Memory
Replicas: 1

Create a Model

General

Name *

Description *

Build

File *



Function *

Example Input ⓘ

```
{  
    "feature": "0, 65, 0, 137, 21.95, 83, 19.42, 111, 9.4, 6, 3.43, 4"  
}
```



Example Output ⓘ

```
{ "result": "value" }
```



Kernel

- Python 2
- Python 3
- R

Comment

Deployment

Engine Profile



Replicas



[Set Environmental Variables](#)

If all parameters are set, you can hit the ‘Deploy Model’ button. Wait till the model is deployed. This will take several minutes.

Models

Model	Status	Replicas	CPU	Memory	Created By	Deployed By	Last Deployed	Actions
My Churn Prediction Model	Building	0 / 1	0	0 GiB	andre	andre	Oct 12, 2018, 9:13 AM	<button>Stop</button>

Step 3: Test the deployed model

After the several minutes, your model should get to the ‘Deployed’ state.

Models

Model	Status	Replicas	CPU	Memory	Created By	Deployed By	Last Deployed	Actions
My Churn Prediction Model	Deployed	1 / 1	1	2 GiB	andre	andre	Oct 12, 2018, 9:15 AM	<button>Stop</button>

Now, click on the Model Name link, to go to the Model Overview page. From the that page, hit the ‘Test’ button to check if the model is working.

My Churn Prediction Model

Overview	Deployments	Builds	Monitoring	Settings
Description	My Churn Prediction Model	Shell	Python	R
Sample Code	<pre>curl -H "Content-Type: application/json" -X POST http://cdsw.52.211.207.216.nip.io/api/altus-ds-1/models/call-model -d '{"accessKey":"me4h5i8qvug7ywd3xbw21ew60wj1fx","request":{"feature":"0, 65, 0, 137, 21.95, 83, 19, 42, 111, 9.4, 6, 3.43, 4"}}'</pre>			
Test Model	<p>Input</p> <pre>{ "feature": "0, 65, 0, 137, 21.95, 83, 19.42, 111, 9.4 , 6, 3.43, 4" }</pre> <p><button>Test</button> <button>Reset</button></p>			
Model Details	<p>Model Id: 10 Deployment: 49 Build: 10 Deployed By: andre Comment: Initial revision. Kernel: python2 Engine Image: Base Image v5 File: predict_churn_sklearn.py Function: predict</p>			
Model Resources	<p>Replicas: 1 Total CPU: 1 vCPUs Total Memory: 2 GiB</p>			

If your model is working, you should receive an output similar like this:

Test Model

Input

```
{  
  "feature": "0, 65, 0, 137, 21.95, 83, 19.42, 111, 9.4  
, 6, 3.43, 4"  
}
```

[Test](#) [Reset](#)

Result

Status	● success
Response	{ "result": 1 }
Replica ID	my-churn-prediction-model-10-49-599d6548d8-x9cmp

The green color with success is telling that our REST call to the model is technically working. And if you examine the response: {"result": 1}, it returns a 1, which mean that customer with these features is likely to churn.

Now, lets change the input parameters and call the predict function again. Put the following values in the Input field:

```
{  
  "feature": "0, 95, 0, 88, 26.62, 75, 21.05, 115, 8.65, 5, 3.32, 3"  
}
```

Test Model

Input

```
{
  "feature": "0, 95, 0, 88, 26.62, 75, 21.05, 115, 8.65
, 5, 3.32, 3"
}
```

[Test](#) [Reset](#)

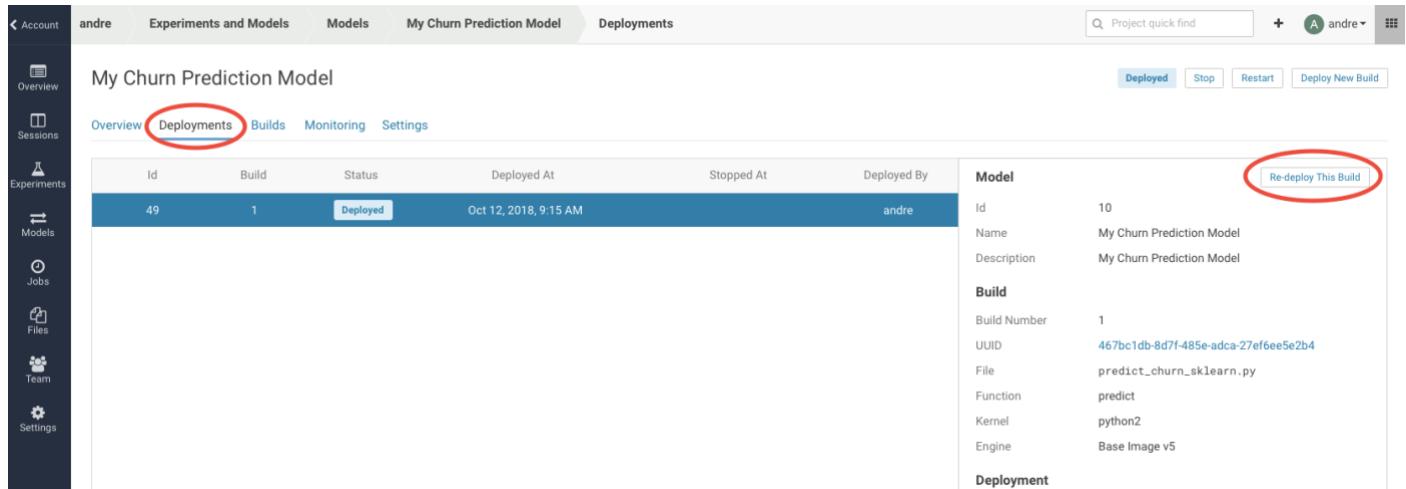
Result

Status	● success
Response	<pre>{ "result": 0 }</pre>
Replica ID	my-churn-prediction-model-10-49-599d6548d8-x9cmp

With these input parameters, the model returns 0, which mean that the customer is not likely to churn.

Step 4: Model Administration

When a model is deployed, Cloudera Data Science Workbench allows you to specify a number of replicas that will be deployed to serve requests. For each active model, you can monitor its replicas by going to the model's Monitoring page. On this page you can track the number of requests being served by each replica, success and failure rates, and their associated stderr and stdout logs. Depending on future resource requirements, you can increase or decrease the number of replicas by re-deploying the model.



The screenshot shows the Cloudera Data Science Workbench interface. The left sidebar has icons for Account, Overview, Sessions, Experiments, Models, Jobs, Files, Team, and Settings. The main navigation bar includes 'Experiments and Models', 'Models', 'My Churn Prediction Model', 'Deployments', and a search bar. Below the navigation is a header for 'My Churn Prediction Model' with tabs for 'Overview', 'Deployments' (which is active), 'Builds', 'Monitoring', and 'Settings'. The 'Deployments' table shows one entry:

ID	Build	Status	Deployed At	Stopped At	Deployed By
49	1	Deployed	Oct 12, 2018, 9:15 AM		andre

To the right of the table is a detailed view of the deployment:

- Model:**
 - ID: 10
 - Name: My Churn Prediction Model
 - Description: My Churn Prediction Model
- Build:**
 - Build Number: 1
 - UUID: 467bc1db-8d7f-485e-adca-27ef6ee5e2b4
 - File: predict_churn_sklearn.py
 - Function: predict
 - Kernel: python2
 - Engine: Base Image v5
- Deployment:** (This section is currently empty)

A red circle highlights the 'Deployments' tab in the navigation bar, and another red circle highlights the 'Re-deploy This Build' button in the deployment details.

Cloudera - Data Science Labs

When you get to the re-deployment page, you can increase the number of replica's.

The screenshot shows the 'Build' and 'Deployment' sections. In the 'Deployment' section, the 'Replicas' input field is highlighted with a red circle and contains the value '1'. Below this field is a 'Deploy Model' button, which is also circled in red.

In order not to overload the cluster, hit the 'Cancel' button to return to the running model page.

Now, navigate to the 'Monitoring' tab.

The screenshot shows the 'Monitoring' tab selected. Below it is a table with columns: Replica, Status, Received, Processed, Success, Failure, Error, Busy, Not Ready, and Restart. One row in the table is highlighted with a blue background. At the bottom of the page, there is a log pane showing several log entries with timestamps and log levels (INFO, DEBUG, etc.). The most recent log entry is highlighted with a red bar.

Several statistics of the model are displayed, like the number of times the model has been called, have been processed, etc.

Logfile information is also available here. The most recent logs are at the top of the pane (see image). stderr logs are displayed next to a red bar while stdout logs are by a green bar. Note that model logs and statistics are only preserved so long as the individual replica is active. When a replica restarts (for example, in case of bad input) the logs also start with a clean slate.

Cloudera - Data Science Labs

Now, navigate to the ‘Settings’ tab.

The screenshot shows the 'My Churn Prediction Model' settings page. On the left is a sidebar with icons for Overview, Sessions, Experiments, Models, Jobs, Files, and Team. The main area has tabs for Overview, Deployments, Builds, Monitoring, and Settings, with 'Settings' circled in red. Below are fields for Name ('My Churn Prediction Model'), Description ('My Churn Prediction Model'), and Access Key ('me4h518qvug7yuwd3xbw21ew60wjlfx'). A 'Regenerate' button is next to the Access Key field. A blue 'Update' button is at the bottom. A 'Danger Zone' section contains a warning: 'Warning! Deleting a model is irreversible. All model builds and deployments history will be deleted.' A red 'Delete Model' button is at the bottom of this section.

On the settings tab, you will be able to find the “Access Key” that is needed in order to call the model with a REST webservice call.

That completes our lab with models. Please, freeup some resources for other people and new projects. So stop from the Models page, also stop your deployed model.

The screenshot shows the 'Models' page. The sidebar includes icons for Overview, Sessions, and Experiments. The main table lists one model: 'My Churn Prediction Model' with status 'Deployed', 1 replica, 1 CPU, 2 GB memory, created by 'andre', deployed by 'andre' on 'Oct 12, 2018, 9:15 AM'. A 'New Model' button is at the top right. A 'Stop' button with a dropdown arrow is circled in red in the 'Actions' column for the listed model.

Model	Status	Replicas	CPU	Memory	Created By	Deployed By	Last Deployed	Actions
My Churn Prediction Model	Deployed	1 / 1	1	2 GB	andre	andre	Oct 12, 2018, 9:15 AM	Stop ▾

Lab 6 - SparklyR

We've focused on python integration, but just to show we can do similar things with R, let's take a look at the R programs and execute them.

This lab requires that R is setup correctly. We provided instructions [earlier](#), but if you skipped them then do this:

Ensure that you've started up an R session (**2vCPU, 4GiB engine**) and executed lines 37 through 42 from the README.md file. You'll only have to do this once for this project so if you've already done it don't do it again.

1. Create a new R Session.

Start New Session

Engine Image - Configure

Base Image v6 - docker.repository.cloudera.com/cdsw/engine:6

Select Engine Kernel

- Python 2
- Python 3
- Scala
- R

Select Engine Profile

1 vCPU / 4 GiB Memory ▾

0 GPUs ▾

Launch Session

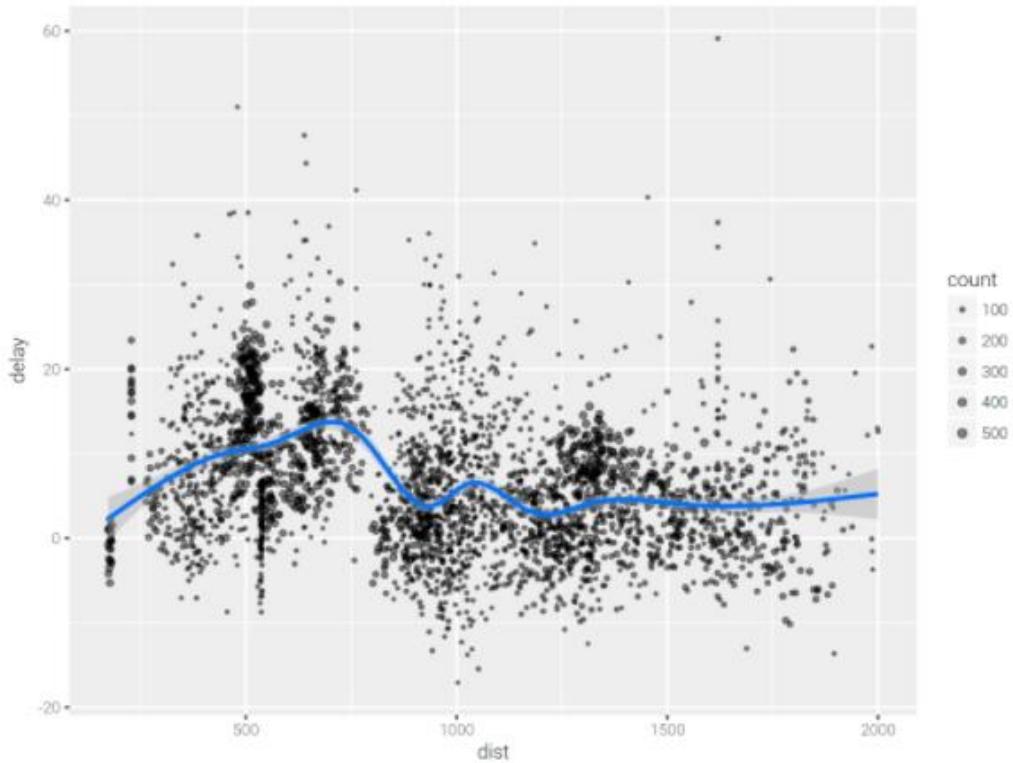
or

Run Experiment..

2. Select (and run) 4_sparklyr.R

The right hand side output window should (eventually) look like this (more or less - depending on your screen real-estate):

```
geom_smooth()` using method = 'gam'
```



Machine Learning

You can orchestrate machine learning algorithms in a Spark cluster via the machine learning functions within sparklyr. connect to a set of high-level APIs built on top of DataFrames that help you create and tune machine learning workflow

In this example we'll use `ml_linear_regression` to fit a linear regression model. We'll use the built-in mtcars dataset, and predict a car's fuel consumption (`mpg`) based on its weight (`wt`) and the number of cylinders the engine contains (`cyl`). In each case that the relationship between `mpg` and each of our features is linear.

```
copy mtcars into spark  
> mtcars_tbl <- copy_to(sc, mtcars, overwrite = TRUE)  
transform our data set, and then partition into 'training', 'test'  
> partitions <- mtcars_tbl %>%  
>
```

Have a read through the Machine Learning section

3. Can you figure out some of the things it's doing? If you know R, and if you know sparklyr, then you can get detailed; if you don't know R then simply 'collapse' the output and see if you can make sense of the analysis without looking at any code ... hopefully you can!

Lab 7 - Shiny

R has a great interactive experience using the shiny package. In this lab we'll create an interactive histogram and you can work with it to find out the frequency distribution of the period between Yellowstone Geyser eruptions!

1. clear the R session screen and then run the 5_shiny.R application

This will start up a shiny server in the local context (line 16), and connect it to the server/ui code (in server.R and ui.R, respectively).

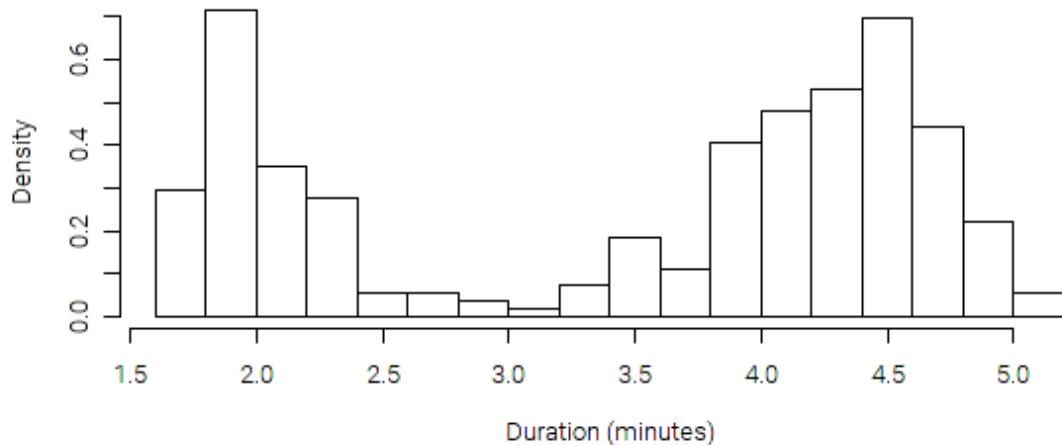
Your users are then presented with a nice little application where they can experiment with changing some of the parameters and graphing features of R's base graphics histogram plot!

Number of bins in histogram (approximate):

20

- Show individual observations
- Show density estimate

Geyser eruption duration



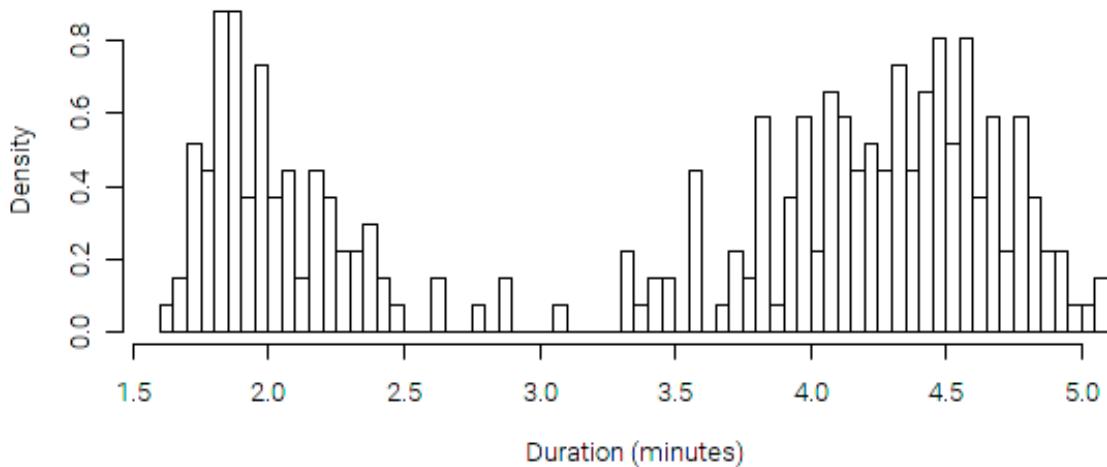
2. You can also change the number of bins in histogram to 50

Number of bins in histogram (approximate):

50

- Show individual observations
- Show density estimate

Geyser eruption duration



Try generating a share link and opening up the share in another browser window - amazingly enough each browser share is independent, allowing your users to share the same underlying experience, but with their individual data inputs!

When you've finished here then stop the R session just to free up some resources.

Have you Stopped the R Session ?

Question: Is this kind of interactivity with data likely to change the way your business users understand and appreciate the work of the Data Scientists?

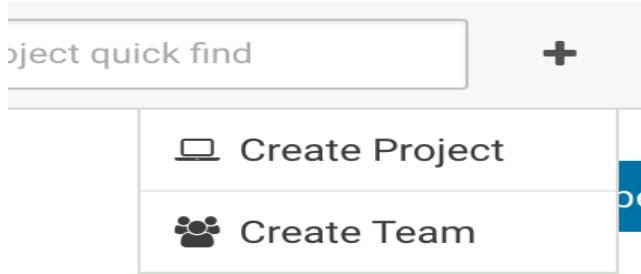
Lab 8 - Scala

In this lab we show how you can use the ‘Template’ mechanism to get started with a simple Scala example. Note that the built in templates and example code aren’t written with multiple users in mind, so you might see file access and permission errors due to the fact that other students might’ve created or deleted files before you!:

1. Navigate to the project space by selecting “project”:



2. Create a new project by hitting the ‘+’ button on the top right and selecting ‘create project’:



3. In the Create new Project window that comes up provide a name for your new project ('Scala', for example), and then choose the Scala template in the Initial Setup drop down menu:

Create a New Project

Project Name

Scala

Project Visibility

- Private** - Only added collaborators can view the project.
- Public** - All authenticated users can view this project.

Initial Setup

Blank

Template

Local

Git

Scala

Templates include example code to help you get started.

Create Project

4. Create the project. You'll see the File Browser view onto the project:

The screenshot shows the Cloudera Data Science Labs interface. On the left is a sidebar with icons for Account, Overview, Sessions, Experiments, Models, Jobs, Files, Team, and Settings. The main area is titled 'Scala' and contains sections for 'Models' (with a note about no models yet), 'Jobs' (with a note about no jobs yet), and 'Files'. The 'Files' section lists several Scala files and configuration files, including 'auction-analysis.scala', 'log4j.properties', 'pi.scala', 'README.md', 'spark-defaults.conf', and 'wordcount.scala'. A download button and a 'New' button are at the top right of the file list.

5. Hit 'Open Workbench' in the top right and let's go run some Scala code:
6. Start a Scala session

Select Engine Kernel

- Python 2
 Python 3
 Scala
 R

Select Engine Profile

1 vCPU / 2 GiB Memory

Launch Session

or Run Experiment..

7. The Scala example project includes its own data set that needs to be moved into HDFS, since that is where the scala code expects to find it. Open a terminal and execute the following shell commands to do this. Note how you have ready access to your own project's data (purely local to you) and the secure (and massive) HDFS cluster:

8. Open Terminal

9. `hdfs dfs -put -f data /user/$HADOOP_USER_NAME`

```
cdsw@44bi8vk24giecnf6:~$ hdfs dfs -put -f data /user/$HADOOP_USER_NAME
cdsw@44bi8vk24giecnf6:~$
```

10. Open example -> kmeans.scala, and edit the file as follows:

Cloudera - Data Science Labs

Remove line 8, by adding // in front. This will mark it as comments and not execute it.

```
7 //load local data to hdfs  
8 // "hdfs dfs -put data/kmeans_data.txt /tmp" !  
9
```

Modify line 11. The location of the input files is in your training directory. Make sure you refer to your training user number.

```
10 //example kmeans clustering script
11 val data = sc.textFile("/user/training20/kmeans_data.txt")
```

11. When the file is modified, run the program.

```
examples/kmeans.scala
File Edit View Navigate Run ➤ examples/kmeans.scala

Scala ⚡
auction-analysis.scala
+ data
examples
  kmeans.scala
  movieALS.scala
log4j.properties
pi.scala
README.md
spark-defaults.conf
wordcount.scala

1 import org.apache.spark.mllib.clustering.{KMeans, KMeansModel}
2 import org.apache.spark.mllib.linalg.Vectors
3 import org.apache.commons.io.FileUtils
4 import java.io.File
5 import sys.process._
6
7 //Load local data to hdfs
8 //hdfs dfs -put data/kmeans_data.txt /tmp/
9
10 //example kmeans clustering script
11 val data = sc.textFile("/user/training20/kmeans_data.txt")
12 val parsedData = data.map(s => Vectors.dense(s.split(' ').map(_.toDouble))).cache()
13
14 // Cluster the training data set into two classes with KMeans
15 val numClusters = 2
16 val numIterations = 20
17 val clusters = KMeans.train(parsedData, numClusters, numIterations)
18
19 // Evaluate clustering by computing Within Set Sum of Squared Errors
20 val WSSSE = clusters.computeCost(parsedData)
21 println(s"Within Set Sum of Squared Errors = $WSSSE")
22
23 // Save the model
24 val output = "output/KMeansExample/KMeansModel"
25 FileUtils.deleteQuietly(new File(output))
26 clusters.save(sc, output)
27
28 //example of loading and predicting on the model we created
29 val sameModel = KMeansModel.load(sc, output)
30 sameModel.clusterCenters.zipWithIndex.foreach { case (center, idx) =>
31   println(s"Cluster Center $idx: ${center}")
32 }
33 sameModel.predict(Vectors.dense(7.5, 6))
34

> import java.io.File
> import sys.process._
load local data to hdfs
"hdfs dfs -put data/kmeans_data.txt /tmp" !

example kmeans clustering script

> val data = sc.textFile("/user/training20/kmeans_data.txt")
> val parsedData = data.map(s => Vectors.dense(s.split(' ').map(_.toDouble))).cache()
Cluster the training data set into two classes with KMeans

> val numClusters = 2
> val numIterations = 20
> val clusters = KMeans.train(parsedData, numClusters, numIterations)
Evaluate clustering by computing Within Set Sum of Squared Errors

> val WSSSE = clusters.computeCost(parsedData)
> println(s"Within Set Sum of Squared Errors = $WSSSE")
Within Set Sum of Squared Errors = 8.1199999999994547

Save the model

> val output = "output/KMeansExample/KMeansModel"
> FileUtils.deleteQuietly(new File(output))
| false

> clusters.save(sc, output)
example of loading and predicting on the model we created

> val sameModel = KMeansModel.load(sc, output)
> sameModel.clusterCenters.zipWithIndex.foreach { case (center, idx) =>
    println(s"Cluster Center $idx: ${center}")
  }
Cluster Center 0: [9.1,9.1,9.1]
Cluster Center 1: [0.1,0.1,0.1]

> sameModel.predict(Vectors.dense(7.5, 6))
| 0


```

Question: How will you use templates when demonstrating CDSW to your friends and colleagues?

Remember to stop your scala Session

Lab 9 - Scheduling Jobs

It's often the case that you need to execute tasks on a periodic basis, and to execute one or more tasks once some other task has succeeded. Obviously there are sophisticated workflow engines but for simple workflows CDSW has a jobs system built in.

This lab goes through the mechanics of creating a simple multi-step job process.

1. Open up your 'cdsw workshop' project to get to this screen:

Name	Size	Last Modified
-	-	1 hour ago
data	3.00 kB	1 hour ago
R	1.60 kB	1 hour ago
spark-warehouse	3.39 kB	1 hour ago
1_python.py	2.50 kB	1 hour ago
2_pyspark.py	769 B	1 hour ago
3_tensorflow.py	1.74 kB	1 hour ago
4_sparklyr.R	536 B	1 hour ago
5_shiny.R	74 B	1 hour ago
README.md	689 B	1 hour ago
server.R	1.09 kB	1 hour ago
spark-defaults.conf	1.48 kB	27 minutes ago
ui.R	-	-
utils.py	-	-
utils.pyc	-	-

2. You need to be in a project to create a Job
3. Select the 'new job' link in the middle of the page, and you'll get to the following screen (there are other ways of getting to this next screen - its an exercise for the student to figure out what they might be):

Create a Job

General

Name

Script
 

Engine Kernel
 Python 2
 Python 3
 Scala
 R

Schedule

Engine Profile

GPUs

Timeout In Minutes (optional) Kill on Timeout

.Jobs exceeding timeout send warning email if notifications enabled.

4. Create a job that will be triggered manually and will execute the 1_python.py. Here are the parameters to do that:

Name	My New Job
Script	1_python.py
Engine Kernel	Python 3

5. Leave everything else as **default**. Scroll down and hit 'Create Job'. You should get to this screen:

Cloudera - Data Science Labs

test CDSW Demo Test Jobs

Project quick find + test

Overview Sessions Experiments Models Jobs Files Team Settings

License Expires in 58 days

1.4.0.431664 (01ac70a)

6. Here you can see that you have a job ('My New Job'). It's never been run, and it has no dependencies.
7. Let's make other jobs depend on this one: Click the '+ Add Job Dependency' grayed out button and add a new job that has a dependency on 'My New Job'. The parameters are:

Name	Job 2
Script	2_pyspark.py
Engine Kernel	Python 3

Name
Job 2

Script
2_pyspark.py 

Engine Kernel
 Python 2
 Python 3
 Scala
 R

Schedule
Dependent 
My New Job 

Engine Profile
0.5 vCPU / 2 GiB Memory 

GPUs
0 GPUs 

Timeout In Minutes (optional) Kill on Timeout
Jobs exceeding timeout send warning email if notifications enabled.

[Set Environmental Variables](#)

8. Scroll down and 'Create Job'. You'll now see a page like this:

Cloudera - Data Science Labs

The screenshot shows the Cloudera Data Science Workshops interface. The top navigation bar includes links for Chrome, File, Edit, View, History, Bookmarks, People, Window, Help, and tabs for EC2 Management Console, CDSW-workshop-2 Environment, Python 2 session (Base Image v1), R session (Base Image v1), Jobs - admin / cdsw workshop, and New Tab. The user is logged in as Toby Ferguson. The left sidebar has sections for Account (admin), Overview, Jobs (selected), Sessions, Files, Team, and Settings. A license status message at the bottom left says "License Expires in 59 days". The main content area is titled "Jobs" and shows "Job Dependencies for My New Job". It displays two jobs: "My New Job" and "Job 2". A dependency graph shows "My New Job" pointing to "Job 2". Below the graph is a table of jobs:

Name	Runs / Failures	Duration	Status	Latest Run	Actions
Job 2	0 / 0	00:00	Not Yet Run	-	<button>Run</button>
My New Job	0 / 0	00:00	Not Yet Run	-	<button>Run</button>

9. So here we can see that 'Job 2' depends upon 'My New Job' (although you can run each manually, if you so choose).
10. Lets add another job that will run in parallel with Job2:
11. Click 'New Job' in the top right corner and create another job that depends upon 'My New Job'. The parameters you'll need are:

Name	R Job
Script	4_sparklyr.R
Engine Kernel	R
Schedule	Dependent / My New Job

General

Name

Script
 

Engine Kernel
 Python 2
 Python 3
 Scala
 R

Schedule
 
 

Engine Profile
 

Timeout In Minutes (optional) Kill on Timeout

Jobs exceeding timeout send warning email if notifications enabled.

[Set Environmental Variables](#)

12. Create the job and you'll see this:

Cloudera - Data Science Labs

The screenshot shows the Cloudera Data Science Workshops interface. On the left is a sidebar with icons for Overview, Jobs, Sessions, Files, Team, and Settings. The main area has tabs for Overview, Jobs, Sessions, Files, and Team. The 'Jobs' tab is selected. A sub-header 'Job Dependencies for My New Job' is shown. Below it is a graph with three nodes: 'My New Job' (top), 'Job 2' (middle), and 'R job' (bottom). Arrows point from 'My New Job' to both 'Job 2' and 'R job'. To the right is a table of jobs:

Name	Runs / Failures	Duration	Status	Latest Run	Actions
R job	0 / 0	00:00	Not Yet Run	-	<button>Run</button>
Job 2	0 / 0	00:00	Not Yet Run	-	<button>Run</button>
My New Job	0 / 0	00:00	Not Yet Run	-	<button>Run</button>

A license status message 'License Expires in 59 days' is at the bottom left, and a build number '1.0.1 (052787a)' is at the bottom right.

13. Lets run it all - hit the 'Run' button next to 'My New Job' (bottom of the list of jobs). You should see the job get scheduled, run, complete, and then the next two jobs should likewise get scheduled, run and complete:

The screenshot shows the same interface after running 'My New Job'. The 'Status' column for 'My New Job' now shows 'Success' instead of 'Not Yet Run'. The 'Latest Run' column shows 'just now'. The other two jobs ('Job 2' and 'R job') have also been run and completed successfully, as indicated by their 'Success' status and 'just now' latest run times.

Question: How will a job scheduler reduce the effort required for you to build simple pipelines?

Question: What other facilities surrounding a job did we not explain? What do you think those other parameters might do?

Challenge Lab 10 – Experiments

Starting with version 1.4, Cloudera Data Science Workbench allows data scientists to run batch experiments that track different versions of code, input parameters, and output (both metrics and files).

Challenge

As data scientists iteratively develop models, they often experiment with datasets, features, libraries, algorithms, and parameters. Even small changes can significantly impact the resulting model. This means data scientists need the ability to iterate and repeat similar experiments in parallel and on demand, as they rely on differences in output and scores to tune parameters until they obtain the best fit for the problem at hand. Such a training workflow requires versioning of the file system, input parameters, and output of each training run.

Without versioned experiments you would need intense process rigor to consistently track training artifacts (data, parameters, code, etc.), and even then it might be impossible to reproduce and explain a given result. This can lead to wasted time/effort during collaboration, not to mention the compliance risks introduced.

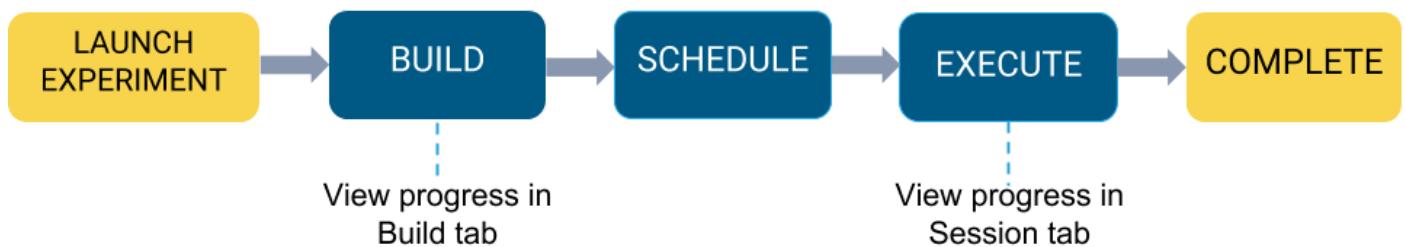
Solution

Starting with version 1.4, Cloudera Data Science Workbench uses experiments to facilitate ad-hoc batch execution and model training. Experiments are batch executed workloads where the code, input parameters, and output artifacts are versioned. This feature also provides a lightweight ability to track output data, including files, metrics, and metadata for comparison.

Concepts

The term experiment refers to a non interactive batch execution script that is versioned across input parameters, project files, and output. Batch experiments are associated with a specific project (much like sessions or jobs) and have no notion of scheduling; they run at creation time. To support versioning of the project files and retain run-level artifacts and metadata, each experiment is executed in an isolated container.

Lifecycle of an Experiment



Step 1: Create a new project

Go to the homepage of your Data Science workbench, and create a 'New' project.

Call the new repository something like Experiments and Models.

Create the repository as a clone of the github repository:

https://github.com/andremolenaar/experiments_models_tellarius.git

Create a New Project

Project Name

Experiments and Models

Project Visibility

- Private** - Only added collaborators can view the project.
- Public** - All authenticated users can view this project.

Initial Setup

Blank

Template

Local

Git

<https://github.com/andremolenaar/dsfortelcoCDSW.git>

Create Project

Start a workbench with a Python 3 and 4 GB of memory.

When the workbench is available, open a terminal window and make the cdsw-build.sh program executable. Use the following command to do that:

```
chmod +x setup.sh  
./setup.sh
```

The screenshot shows a terminal window titled "Untitled Session" running on the Cloudera Data Science Workbench. The terminal has a red circle highlighting the "Terminal access" button in the top bar. The terminal output shows the following commands and their results:

```
https://tty-l3qapsp5u0d1nyp4.cdsweb.cloudera.tellarius.eu/8b8i1gl84yt6j55s/  
Project workspace: /home/cdsw  
Kerberos principal: training20@CLOUDERA.TELLARIUS.EU  
Runtimes:  
  R: R version 3.4.1 (--) -- "Single Candle"  
  Python 2: Python 2.7.11  
  Python 3: Python 3.6.1  
  Java: java version "1.8.0_131"  
Git origin: https://github.com/andremolenaar/experiments_models_tellarius.git  
cdsw@l3qapsp5u0d1nyp4:~$ chmod +x setup.sh  
cdsw@l3qapsp5u0d1nyp4:~$ ./setup.sh  
cdsw@l3qapsp5u0d1nyp4:~$
```

Step 2: Examin dsfortelco_sklearn_exp.py

Open the file “dsfortelco_sklearn_exp.py”. This is a python program that builds a churn model to predict customer churn (the likelihood that this customer is going to stop his subscription with his telecom operator). There is a dataset available on hdfs (/user/trainingXX/churn_all.csv), with customer data, including a churn indicator field.

The program is going to build a churn prediction model using the Random Forest algorithm. Random forests are ensembles of decision trees. Random forests are one of the most successful machine learning models for classification and regression. They combine many decision trees in order to reduce the risk of overfitting. Like decision trees, random forests handle categorical features, extend to the multiclass classification setting, do not require feature scaling, and are able to capture non-linearities and feature interactions.

spark.mllib supports random forests for binary and multiclass classification and for regression, using both continuous and categorical features. spark.mllib implements random forests using the existing decision tree implementation. Please see the decision tree guide for more information on trees.

The Random Forest algorithm expects a couple of parameters:

- numTrees: Number of trees in the forest.
Increasing the number of trees will decrease the variance in predictions, improving the model’s test-time accuracy.
Training time increases roughly linearly in the number of trees.
- maxDepth: Maximum depth of each tree in the forest.
Increasing the depth makes the model more expressive and powerful. However, deep trees take longer to train and are also more prone to overfitting.
In general, it is acceptable to train deeper trees when using random forests than when using a single decision tree. One tree is more likely to overfit than a random forest (because of the variance reduction from averaging multiple trees in the forest).

In the dsfortelco_pyspark_exp.py program, these parameters can be passed to the program at runtime. In the lines 38 and 39, these parameters are passed to python variables:

```
param_numTrees=int(sys.argv[1])  
param_maxDepth=int(sys.argv[2])
```

Also note that at the lines 69 and 70, the quality indicator for the Random Forest model, are written back to the Data Science Workbench repository:

```
cdsw.track_metric("auroc", auroc)  
cdsw.track_metric("ap", ap)
```

These indicators will show up later in the Experiments dashboard.

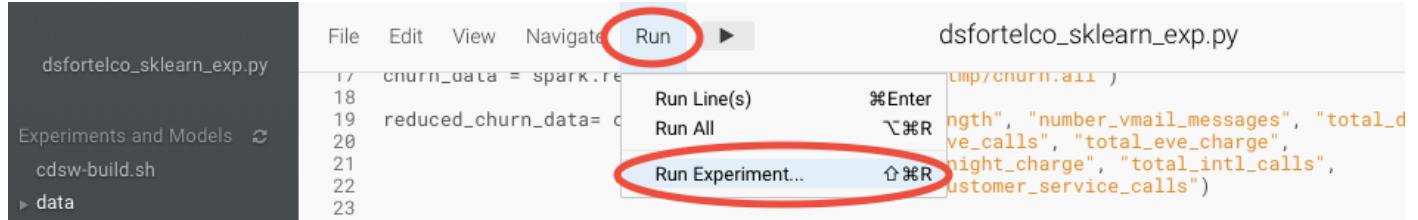
Step 3: Run the experiment for the first time

Now, run the experiment using the following parameters:

numTrees = 40

numDepth = 20

From the menu, select Run -> Experiments.



Specify the arguments for this run, by typing the numbers behind the arguments field. Note that these fields are separated by a space and that there is no comma (,)

Run New Experiment

Script	dsfortelco_sklearn_exp.py
Arguments	40 20
Engine Kernel	<input checked="" type="radio"/> Python 3 <input type="radio"/> Python 2 <input type="radio"/> Scala <input type="radio"/> R
Engine Profile	0.5 vCPU / 1 GiB Memory
Comment	First Random Forrest Experiment
Cancel Start Run	

Now, in the background, the Data Science Workbench environment will spin up a new docker container, where this program will run.

Step 4: Check the results for the first experiment

Go back to the ‘Projects’ page in CDSW, and hit the ‘Experiments’ button.



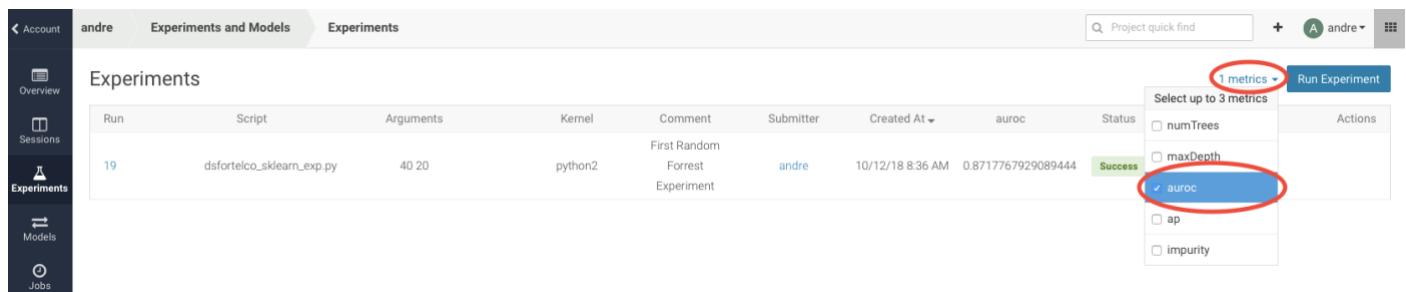
The screenshot shows the CDSW interface with the 'Experiments' tab selected in the sidebar. The main table displays one experiment run:

Run	Script	Arguments	Kernel	Comment	Submitter	Created At	Status	Duration	Actions
19	dsfortelco_sklearn_exp.py	40 20	python2	First Random Forrest Experiment	andre	10/12/18 8:36 AM	Building		<button>Stop</button>

If the Status indicates ‘Running’, you have to wait till the run is completed.

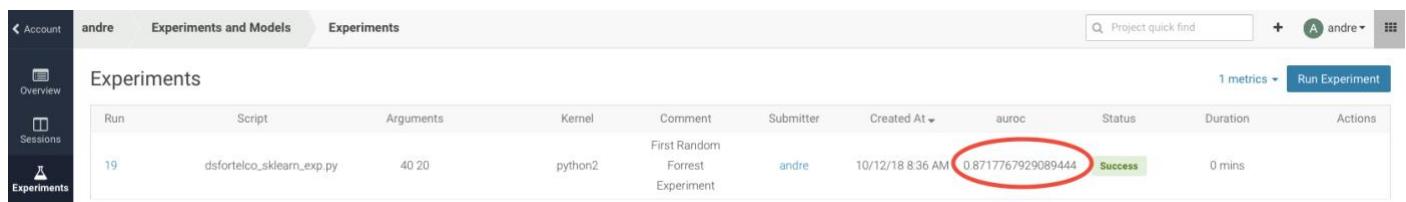
In case the status is ‘Build Failed’ or ‘Failed’, check the log information. This is accessible by clicking on the run number of your experiments. There you can find the session log, as well as the build information.

In case your status indicates ‘Success’, you should be able to see the auroc (Area Under the Curve) model quality indicator. It might be that this value is hidden by the CDSW user interface. in that case, click on the ‘3 metrics’ links, and select the auroc field. It might be needed to de-select some other fields, since the interface can only show 3 metrics at the same time.



The screenshot shows the CDSW interface with the 'Experiments' tab selected in the sidebar. The main table displays one experiment run. A dropdown menu titled 'Select up to 3 metrics' is open, showing several options. The 'auroc' option is highlighted with a blue oval and has a checkmark next to it, indicating it is selected.

When the auroc metric is selected, you will be able to see the value.



The screenshot shows the CDSW interface with the 'Experiments' tab selected in the sidebar. The main table displays one experiment run. The 'auroc' column now shows the value '0.8717767929089444', which is circled in red. The status is shown as 'Success'.

In this example, 0.871

Not bad, but maybe there are better hyper parameter values available.

Step 5: Re run the experiment several times

Now, re-run the experiment 3 more times and try different values for NumTrees and NumDepth. Try the following values:

NumTrees	NumDepth
15	25
25	20
Try something yourself	Try something yourself

When all runs have completed successfully, check which parameters had the best quality (best predictive value). This is represented by the highest 'area under the curve', auroc metric.

Run	Script	Arguments	Kernel	Comment	Submitter	Created At	numTrees	maxDepth	auroc	Status	Duration	Actions
22	dsfortelco_sklearn_exp.py	20 30	python2	Fourth Experiment	andre	10/12/18 8:47 AM	20	30	0.8392243291195386	Success	0 mins	
21	dsfortelco_sklearn_exp.py	25 20	python2	Third Experiment	andre	10/12/18 8:46 AM	25	20	0.8502144506763445	Success	0 mins	
20	dsfortelco_sklearn_exp.py	15 25	python2	Second Experiment	andre	10/12/18 8:46 AM	15	25	0.8178173020922	Success	1 mins	
19	dsfortelco_sklearn_exp.py	40 20	python2	First Random Forrest Experiment	andre	10/12/18 8:36 AM	40	20	0.8717767929089444	Success	0 mins	

In this example, run 21 had the highest auroc value, so that is the model that you would want to use for your business.

Step 6: Save the best model to your environment

Select the run number with the best predictive value, in this example, run number 21.

Run	Script	Arguments	Kernel	Comment	Submitter	Created At	numTrees	maxDepth	auroc	Status	Duration	Actions
22	dsfortelco_sklearn_exp.py	20 30	python2	Fourth Experiment	andre	10/12/18 8:47 AM	20	30	0.8392243291195386	Success	0 mins	
21	dsfortelco_sklearn_exp.py	25 20	python2	Third Experiment	andre	10/12/18 8:46 AM	25	20	0.8502144506763445	Success	0 mins	
20	dsfortelco_sklearn_exp.py	15 25	python2	Second Experiment	andre	10/12/18 8:46 AM	15	25	0.8178173020922	Success	1 mins	
19	dsfortelco_sklearn_exp.py	40 20	python2	First Random Forrest Experiment	andre	10/12/18 8:36 AM	40	20	0.8717767929089444	Success	0 mins	

In the Overview screen of the experiment, you can see that the model in spark format, is captured in the

file 'sklearn_rf.pkl'. Select this file and hit the 'Add to Project' button. This will copy the model to your project directory.

The screenshot shows the 'Experiments and Models' section of the Cloudera Data Science Labs interface. On the left, there's a sidebar with icons for Overview, Sessions, Experiments, Models, Jobs, Files, Team, and Settings. The main area shows a run named 'Run-21' with tabs for Overview, Session, and Build. Under Configuration, it lists a Script as 'dsfortelco_sklearn_exp.py', Arguments as '25 20', a Comment as 'Third Experiment', a Build Snapshot as '77e3b0ce1cf6b010c840e282aeda6d7ffd8dc248', and a Submitter as 'andre'. Under Metrics, it lists numTrees as 25, maxDepth as 20, impurity as gini, auroc as 0.8502144506763445, and ap as 0.680864067072068. To the right, under Output, there's a checked checkbox for 'sklearn_rf.pkl' and a blue-bordered button labeled 'Add to Project'.

Configuration	
Script	dsfortelco_sklearn_exp.py
Arguments	25 20
Comment	Third Experiment
Build Snapshot	77e3b0ce1cf6b010c840e282aeda6d7ffd8dc248
Created At	10/12/18 8:46 AM
Submitter	andre

Metrics	
numTrees	25
maxDepth	20
impurity	gini
auroc	0.8502144506763445
ap	0.680864067072068

Output

sklearn_rf.pkl

Add to Project

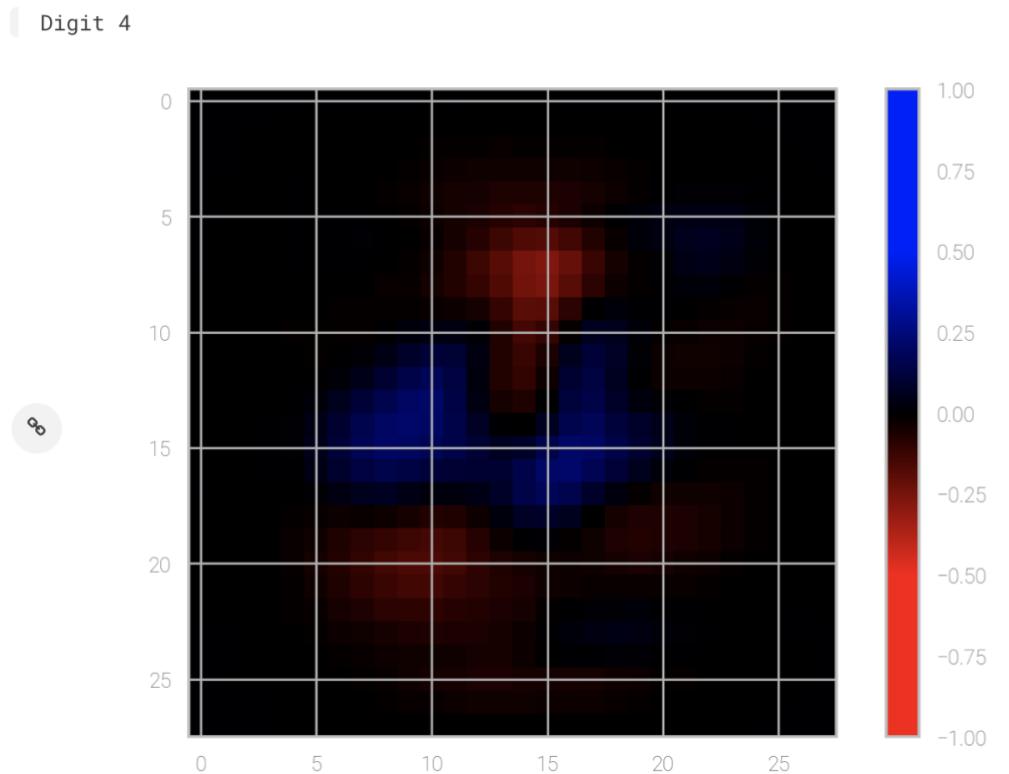
Challenge Lab 11 - Pushing the Boundaries

So you've just heard that the latest thing from Google is [Tensorflow](#) and you're keen to get started. You're going to need to install some customer packages and, more importantly, connect a program providing a web interface so that you can view the results. Your IT department isn't going to help you - you're going to want to do this experiment on your own.

Fortunately CDSW enables you to install custom libraries. This cluster might be managed by IT but you can still get your libraries in there to do the work you need ...

We've done the hard work for you - take a look:

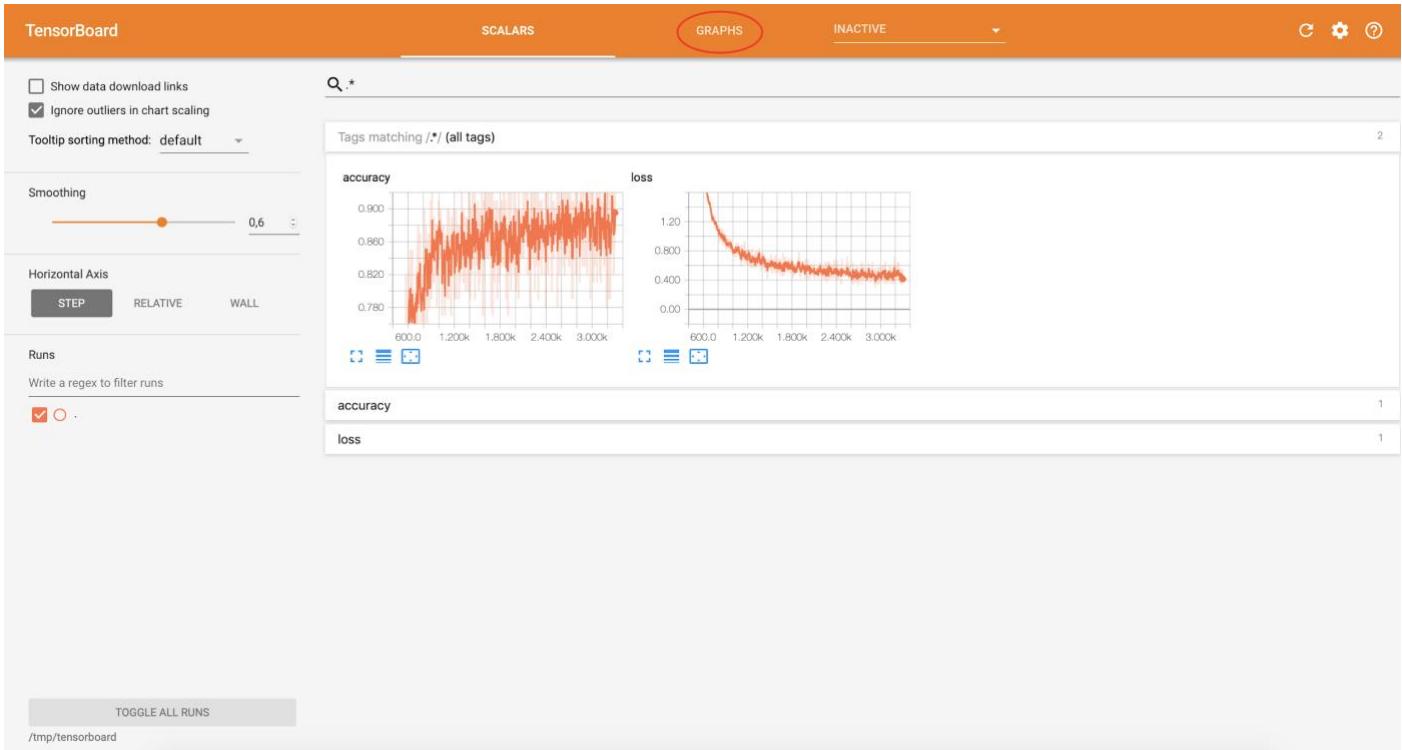
1. Select 3_tensorflow.py
2. Lines 2 through 6 show the imports of the various libraries (we installed them in Step 3)
3. Run 3_tensorflow.py in your current Python session (you might like to 'clear' your output screen first)- the input handwritten number images are shown, along with some images of feature maps for particular numbers



Explore using Tensorboard

```
> utils.start_tensorboard(logs_path, iframe=False)
| Starting Tensorboard at https://ck44c57bnm5xvzq2.cdsweb.cloudera.tellarius.eu...
| Open Tensorboard
| TensorBoard 0.1.8 at http://ck44c57bnm5xvzq2:8080 (Press CTRL+C to quit)
```

When the program has completed, hit the generated ‘Open Tensorboard’ link. This will bring you to the interactive Tensorflow UI, that will help you understand the model.



Have a look at the ‘Graph’ of the model. And take some time to play around with the Tensorboard. When you are ready, close the Tensorboard tab and go back to the Cloudera Data Science Workbench.

Key takeaways: You were able to install and use custom third party libraries, as well as run an application that you can connect to from an external application.

STOP all your Python sessions now (you should only have one but sometimes people get carried away). This will help ensure there are plenty of resources for you and the others in the workshop.

Appendix

Cloudera Documentation

<http://www.cloudera.com/documentation.html>

Cloudera Data Science Workbench

<https://www.cloudera.com/products/data-science-and-engineering/data-science-workbench.html>

CDSW User Guide

https://www.cloudera.com/documentation/data-science-workbench/latest/topics/cdsw_user_guide.html

Troubleshooting Guide

https://www.cloudera.com/documentation/data-science-workbench/latest/topics/cdsw_troubleshooting.html

Recordings

Part 1 - Introduction

<https://www.cloudera.com/content/dam/www/marketing/resources/webinars/introducing-cloudera-data-science-workbench-part1-recorded-webinar.png.landing.html>

Part 2 – A Visual Dive into machine Learning

<https://www.cloudera.com/content/dam/www/marketing/resources/webinars/part-2-visual-dive-into-machine-learning-and-deep-learning.png.landing.html>

Part 3 – Data Science Models into production from beginner to end

<https://www.cloudera.com/content/dam/www/marketing/resources/webinars/models-in-production-a-look-from-beginning-to-end-part3.png.landing.html>

