# Cloudera Machine Learning Workshop

1)  Open CDP, using the "admin" user within the Test Drive link.

Your link should look something like (remember click the link in your email not the link below)
http://login.trycdp.com/auth/realms/trycdp-trialxx/protocol/saml/clients/samlclient?tn=trialxx_admin@trycdp.com&p=
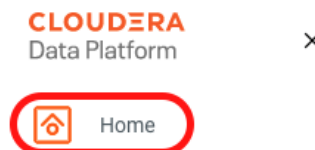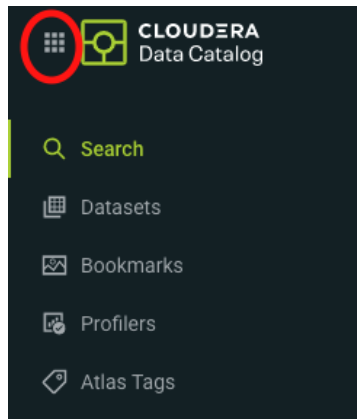X *xx represents the trial user #
*X represents the password

2) Click the "Machine Learning" tile within the CDP Home Screen



How do you get to the CDP Home Screen?
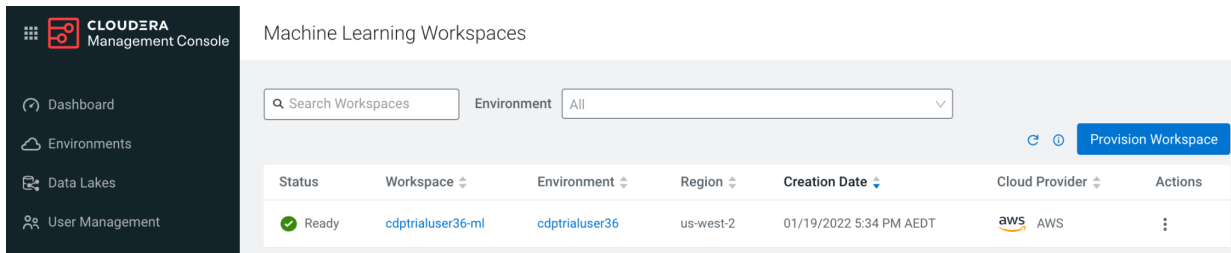
-   From any experience such as "Data Catalog", click the 9 square at the top left and then
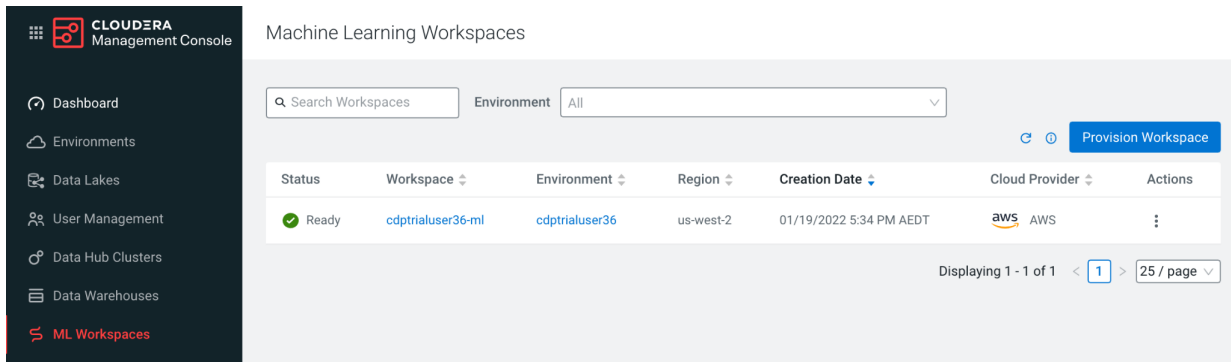    click "Home"

## Part 1 - REQUIRED IN ORDER FOR LAB TO WORK PROPERLY: Create a CML Project

**(A Workspace should already be provisioned for you)**

Workspaces are the heart of the Cloudera Machine Learning (CML)



A Workspace is a small cluster that runs on a kubernetes service to provide teams of data scientists to develop, test, train, and ultimately deploy machine learning models. **Click into the Workspace by clicking the Workspace name.**



You can visualize all of the Projects and Resources that are part of the Projects page. Next we will create a Project where we will develop and deploy models along with other CML features. **Click on "New Project"**

When creating a new project give a Name, Visibility, and initial configuration.

**Project Name: Telco_churn Visibility: Private Initial Setup: Git ->**
**https://github.com/andy-hansen/cml.git**

## Create a New Project

**Project Name**

Telco_churn

**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/andy-hansen/cml.git

Create Project

**Part 2: CML Project Overview**
Overview gives you access to all the features of a CML project. We will only have files copied from the Github repo currently. Initially it is good to start on the management components of a project.

**Collaborators:**

**For our demo we aren't adding additional collaborators.**

You can give access to other users with certain permissions for the encompassing project so teams of users can collaborate together. You can set up Admins, Contributor, Operator, and Viewer permissions.

**Project Settings:**

Taking a look at Project Settings, this is where you can define several options for the current project. You have the ability to define different engines where your code in CML will run. There are project variables that can be defined and used throughout your code. SSH tunnels can also be configured to connect to other services as needed. More details can be found in our docs here.



**-> Change the Runtime/Engine to Legacy Engine -> Click Save Engine**

**Part 3: CML Sessions and Workbench**

Sessions allow you to perform actions such as run R or Python code. They also provide access to an interactive command prompt and terminal. Sessions will be built on a specified Engine Image, which is a docker container that is deployed onto the Workspace. In addition you can specify how many resources are used per session. **From the Overview page click on New Session**



**Session Name: telco_churn_session_1 Editor: Workbench Kernel: Python 3 Engine Image: Default Resource Profile: 1vCPU/2 GiB Memory**

Then select Start Session

## Start A New Session

⚠ This Project is configured to run on Legacy Engines.
You can try the new ML Runtimes by switching over on the Project Settings page.

To learn more visit the documentation.

**Session Name**

telco_churn_session_1

**Engine**

Editor ⓘ

Workbench ⌄

Kernel ⓘ

Python 3 ⌄

**Engine Image -** Configure

Default engine image - docker.repository.cloudera.com/cloudera/cdsw/engine:15-cml-2021.09-2
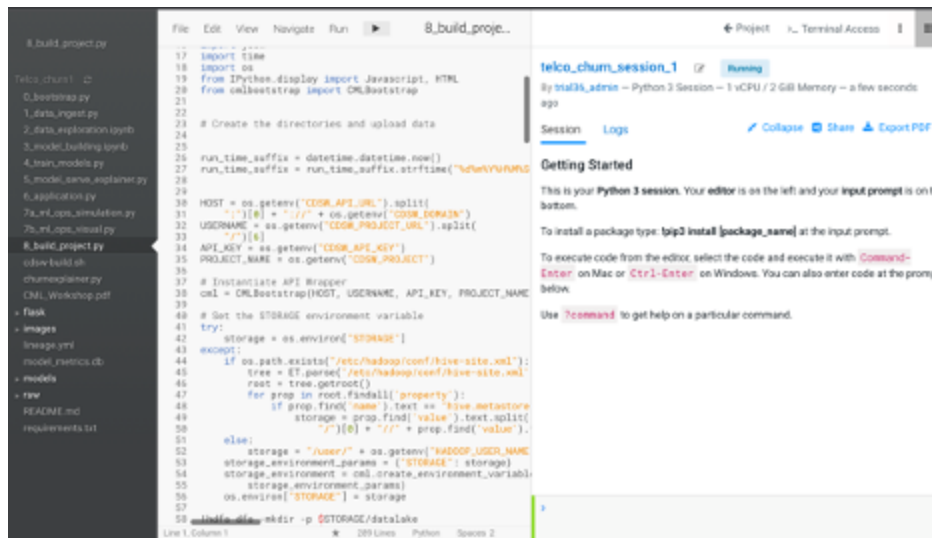
**Resource Profile**

1 vCPU / 2 GiB Memory ⌄

Cancel          **Start Session**

The Workbench is now starting up and deploying a container onto the workspace at this point. Going from left to right you will see the project files, editor pane, and session pane. **Once you see the flashing red line on the bottom of the session pane turn steady green the container has been successfully started.**



**Open up the script 8_build_project.py from the left pane.** In the editor pane we are going to select and run the script in several parts. Throughout the script you will see breaks in the code defined by Part 1, Part 2, Part 3, etc.

Part 1: Will install any required packages to be used. As an example, flask is installed as part of the project. A variable is also set as part of the project. Last but not least we are loading a file from the project as a test dataset and moving that to a s3 location.

**Select and highlight the contents of part 1, then Right click -> Run Line(s)**



Part 2: The telco churn dataset is ingested from s3 and a hive table is created using Spark.
**Select and highlight the contents of part 2, then Right click -> Run Line(s)**


Part 3: Create a CML Job and start the job. A job automates the action of launching an engine, running a script, and tracking the results, all in one batch process. Jobs are created within the purview of a single project and can be configured to run on a recurring schedule. You can customize the engine environment for a job, set up email alerts for successful or failed job runs, and email the output of the job to yourself or a colleague.

**Select and highlight the contents of part 3, then Right click -> Run Line(s)**


Once the Job is started we will look at what was created. **Click the Project button** (top right corner of screen).

**Click on Jobs** to explore the job that was created and details can be found by **Clicking on the Job Name.**

Part 4: Using CML, you can create any function within a script and deploy it to a REST API. In a machine learning project, this will typically be a predict function that will accept an input and return a prediction based on the model's parameters.

**Select and highlight the contents of part 4, then Right click -> Run Line(s)**
When getting to deploying the Model this can take a little time, and is a good spot to stretch the legs and refill your coffee.



Part 5: Applications give data scientists a way to create ML web applications/dashboards and easily share them with other business stakeholders. Applications can range from single visualizations embedded in reports, to rich dashboard solutions such as Tableau. They can be interactive or non-interactive.

Applications stand alongside other existing forms of workloads in CML (sessions, jobs, experiments, models). Like all other workloads, applications must be created within the scope of a project. Each application is launched within its own isolated engine. Additionally, like models, engines launched for applications do not time out automatically. They will run as long as the web application needs to be accessible by any users and must be stopped manually when needed.

**Select and highlight the contents of part 5, then Right click -> Run Line(s)**

When deploying the CML Flask Application you will be provided with a URL to follow at the end of your session pane output. **Click to Open Application UI**



Within the Flask Application UI you can experiment with some of the parameters that will run against the model to predict how likely a customer is to churn:

## Single Prediction View

Churn Probability    0.692

| Field | Value | Score | Options | | | |
|---|---|---|---|---|---|---|
| Contract | Month-to-month | 0.12 | Month-to-month | One year | Two year | |
| Dependents | No | 0 | No | Yes | | |
| DeviceProtection | No | 0 | No | No internet service | Yes | |
| InternetService | Fiber optic | 0.20 | DSL | Fiber optic | No | |
| MonthlyCharges | 97.9 | -0.24 | mean 64.80  min 18.25  max 118.75 | [ ] Submit | | |
| MultipleLines | No | -0.04 | No | No phone service | Yes | |
| OnlineBackup | Yes | 0 | No | No internet service | Yes | |
| OnlineSecurity | No | 0 | No | No internet service | Yes | |
| PaperlessBilling | Yes | 0 | No | Yes | | |
| Partner | No | 0 | No | Yes | | |
| PaymentMethod | Electronic check | 0.04 | Bank transfer (automatic) | Credit card (automatic) | Electronic check | Mailed check |
| PhoneService | Yes | 0 | No | Yes | | |
| SeniorCitizen | No | -0.04 | No | Yes | | |
| StreamingMovies | Yes | 0.08 | No | No internet service | Yes | |
| StreamingTV | Yes | 0.08 | No | No internet service | Yes | |
| TechSupport | Yes | 0 | No | No internet service | Yes | |
| TotalCharges | 315.3 | -0.12 | mean 2283.30  min 18.80  max 8684.80 | [ ] Submit | | |
| gender | Female | 0 | Female | Male | | |
| tenure | 3 | 0.29 | mean 32.42  min 1.00  max 72.00 | [ ] Submit | | |

Part 6: Run the last code snippet in the workbench to complete the project build script (8_build_project.py). This goes through a process of simulating a model that drifts over 1000 calls to the model. The file contains comments with details of how this is done. **Select and highlight the contents of part 6, then Right click -> Run Line(s)**
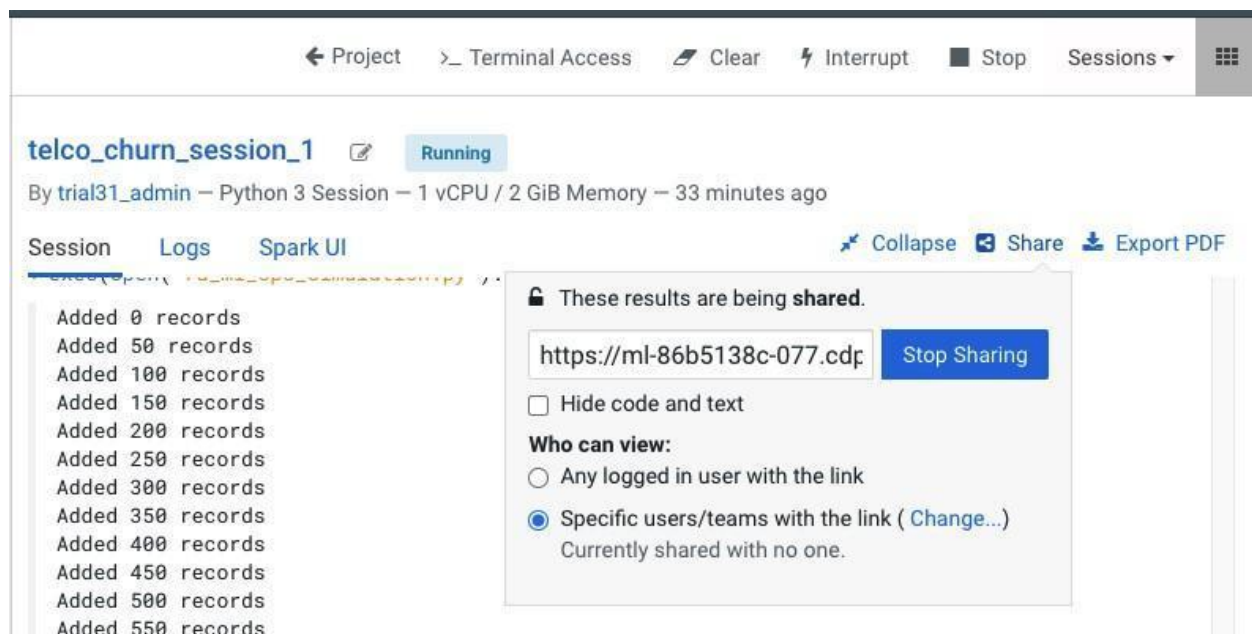


You can also share the workbench session with other users if they would like to view results of the code. A URL can be shared out if users outside of CML would like to view code/results of work.

**Click on share at the top of the Session pane**

A full CML project should now be running with a Job, Model, and Application deployed! Go back to the project and take a look at the Model and Applications page.

Going to the **Model** page will show you the deployed model. Clicking on the Model will give various tabs of monitoring and statistics in addition to previous deployments of the same model.

On the **Applications** page we can see our running flask app we looked at previously.

Models and Applications will continue to run even if the workbench session is stopped or timeout occurs. These will run on engines as part of the CML workspace