

Projet Master 2 Informatique
Apprentissage, Information et Contenu (AIC)

Preparing a starting kit and a challenge bundle

Isabelle Guyon

iguyon@lri.fr

** Please preferably communicate through [Piazza](#) **

Version 4, Sep.23, 2019

Getting started:

[\[Intro slides\]](#)

- (1) Install on your computer **Anaconda Python 3.7**. If you already have another version of Python or Anaconda on your computer, use virtual environments:

```
conda update conda
```

```
conda create -n python37 python=3.7 anaconda
```

Then you will be able to use the 3.7 version by activating the virtual environment:

```
conda activate python37
```

To get out, use:

```
conda deactivate
```

This installation should include python 3.7, the spyder IDE and jupyter-notebook (verify).

- (2) Another possibility will be to use the docker `codalab/codalab-legacy:py3`. More on dockers at the end of this document. If you are going to use OTHER libraries than those provided in Anaconda, definitely use dockers!

- (3) **Clone the starting kit template:** <https://github.com/madclam/m2aic2019/tree/master>. If you are new to Github, take a tutorial. Using Git and Github is essential to your project. This example allows you to build the sample competition <https://codalab.lri.fr/competitions/204>.

- (4) **When done creating your starting kit, push it to your own Github repo (your homework will be your Github repo URL). Turn in your homework on [Piazza](#), tag 2.starting-kit before Sep. 28.**

- (5) Each group has an email grouppname@chalearn.org, which has been used to create an account on: <https://codalab.lri.fr/>.

The **password** is the same for all accounts: **m2_aic++**

1) Starting kit structure:

A starting kit should imperatively have the following structure, do NOT change it:

- **README.md**: Short description of your task and credits.
- **README.ipynb**: Sample code reading data, calling sample code, generating sample submissions.
- **sample_data/**: A SMALL data sub-sample (could be a subset of the training data, but should NOT include any validation or test examples on which the competitors will be tested).
- **sample_result_submission/**: Example prediction files, for the validation and test set.

- **sample_code_submission/**: Self-contained, ready-to-submit example code submission, including the predictive mode class `model.py` and a `metadata` file (with just a comment).
- **ingestion_program/**: Program and libraries needed to run code submissions on the challenge platform and generate prediction results using the `input_data` (labeled training data and unlabeled validation and test data) and `model.py`.
- **scoring_program/**: Program and libraries needed to score the prediction results on the challenge platform using the target values (solution), also called `reference_data`. A library of scoring functions to choose from is provided in `libscore.py`. All you need to do is to put the name of the chosen function in `metric.txt`. If needed, you can supply the code of your own metric in `my_metric.py`. You can also use a scikit-learn metric.
- **html_files**: Documentation files in HTML format.

2) Data preprocessing and formatting:

This year tasks should be limited to **classification** and **regression** problems.

- **Data representation**: For the first (basic) version of the starting kit, all datasets will have to be transformed/preprocessed into a "**feature representation**". This is especially important for text, speech, vision, or any datasets whose raw data consist of images or sequences. For example, for image data, you could use features calculated by a [pre-trained CNN](#), but we do not necessarily recommend it. Those features are really "too good" and "too abstract" to make the problem interesting to the L2 students. A better choice would be to use "classical features" from [OpenCV](#). The groups working on computer vision applications will keep a non-preprocessed version of their dataset, for later (we'll do a second version from raw data for more advanced students).
- **Missing data**: Avoid datasets with a large number of missing values (more than 10%). Represent the missing values with the symbol **NaN**.
- **Binary and categorical variables**: Represent all variables with a **NUMERIC code**. Avoid categorical variables with a large number of values, some of which appear very infrequently. To that end, you may group infrequent categories. Likewise, for binary variables, avoid datasets in which variables are very imbalanced (very sparse data). The target values (solutions) must also be in a numerical format.
- **Number of classes**: Avoid problems with a too large number of classes (>20) and have too few examples per class (less than 1000 in test data).
- **Size dataset**: Do not consider too small datasets of less than 50 000 examples. You must have at least **10 000 test examples** and preferably a large training set and a validation set of the same order of magnitude as the test set. We want enough samples, but the overall preprocessed compressed volume of data should not exceed **300 MB**. Quantizing values to enough precision (but not too much) is recommended, as part of preprocessing. Usually quantizing between 0 and 999 is enough. To further reduce data volume, you may also reduce the number of features by feature selection (this should be the last resort because it is good for the L2 students to struggle with many features). So no limit on the number of features.
- **Data homogeneity**: Avoid data that are non-homogeneous or non-identically-distributed. **Shuffle** the order of your examples randomly.

Data format: All pre-processed data must be produced in [AutoML format](#):

This is imperative because you will all share the same data reader (called by the ingestion program). You will have to split the data into a “**training set**”, a development test set a.k.a. “**validation set**” for the challenge phase 1, and a (final) “**test set**” for the challenge phase 2. You should eventually produce ALL the following files, but, for the first version, produce at least those outlined in boldface:

FileName_train.data # This is the training matrix X

FileName_train.solution # This is the training matrix Y (target values)

FileName_valid.data # This is the first test set called “validation set” used during development

FileName_valid.solution # and its target values (hidden to the participants)

FileName_test.data # This is the FINAL test set

FileName_test.solution # and its target values

FileName_feat.name

FileName_label.name

FileName_feat.type

FileName_public.info

FileName_private.info

We ask that you to:

- + Provide all the requested files (even those not in boldface) because this will greatly help the participants and this be useful documentation for future reference.
- + NOT provide “un-split” data in the simpler form: FileName.data and FileName.solution, because **providing your own data-split** is better for several reasons:

- You have more control and you can, for instance, ensure that each class has a balanced mix of examples (in the case of multiple data sources of multiple types of subjects), or, on the contrary, that you have different types of examples in training and test data (for instance if you want to address speaker independent speech recognition, samples from different speakers should be in training and test data).
- You are sure that the examples are well shuffled.
- You are sure that the example order matches your sample result submission.

Data leakage: This refers to features that inform inadvertently about the target values: running a feature selection algorithm often reveals such features. Typically, there may be a column with sample IDs highly correlated with the target. There are other forms of data leakage. For example if you forget to shuffle the examples and all the examples of the same class appear next to one another. Check https://www.cs.umb.edu/~ding/history/470_670_fall_2011/papers/cs670_Tran_PreferedPaper_LeakinDataMining.pdf for other horror stories.

Sample data and challenge data: You need to prepare 2 datasets in AutoML format:

- **Challenge data:** your actual challenge data. In the challenge platform, it will be split into downloadable **public_data** or **input_data** (everything except the FileName_valid.solution, FileName_test.solution, and FileName_private.info) and **reference_data** not leaving the platform: FileName_valid.solution for phase 1 and FileName_test.solution for phase 2.
- **Sample data:** a small data subset carved out the challenge TRAINING data, for practice purposes only (do not compromise real validation or test data).

3) Sample submissions:

You must prepare sample code submissions put in directory **sample_code_submission/** and **sample_result_submission/**:

1) sample_code_submission/: This directory should contain a Python class `model.py`. This should be a basic but functional baseline method solving your problem. Using the “ingestion program” (see below), you can generate a **sample_result_submission** from your **sample_code_submission**. This directory should also contain a `metadata` file with just a comment. The directory may or may not also contain a pickle with a trained model.

Initially, you can predict all zeros (with the right number of lines).

2) sample_result_submission.zip: This will contain at least 2 files (the last one is optional):

- `dataName_valid.predict`
- `dataName_test.predict`
- `dataName_train.predict`

corresponding to predictions of the **target values** made by your **sample code** for the validation set, the test set, and optionally the training set. There should be one sample per line. For the HADACA project, your participants must predict 2 matrices A and T of dim (N, C) and (C, M). To respect the format, please write matrix A and T' one after the other to get a file with N+M lines and C columns.

4) Ingestion program:

A sample ingestion program is provided (**it should not need to be modified if you respect the AutoML format and the submission format**). The ingestion program receives the code submission of the participants on the server and runs it. However, you can also run it locally to produce a result submission: `python ingestion_program/ingestion.py big_data sample_result_submission ingestion_program sample_code_submission`

where `big_data/` should contain your real challenge data. For debug purposes, you can use `sample_dat/` instead of `big_data/`. **Remove the pickle in sample_code_submission, if any** (may cause an error).

5) Scoring program:

A scoring program is provided. **It should not need to be modified**, however, the scoring metric can be chosen by the organizers (you). Place your code in `my_metric.py`, and indicate the name of your scoring function in `metric.txt`. Instead of writing your own code, you may also select one of the metrics for `libscore.py` and specify it in `metric.txt`:

Binary classification	Multi-class classification	Multi-label classification	Regression
<code>bac_binary</code>	<code>bac_multiclass</code>	<code>bac_multilabel</code>	<code>abs_regression</code>
<code>auc_binary</code>		<code>auc_multilabel</code>	<code>r2_regression</code>
<code>pac_binary</code>	<code>pac_multiclass</code>	<code>pac_multilabel</code>	
<code>f1_binary</code>		<code>f1_multilabel</code>	

Test your metric with:

`python scoring_program/score.py reference_data sample_result_submission scoring_output`

where `reference_data/` should contain your validation and test target values/solutions:

`DataName_valid.solution` for phase 1 and `DataName_test.solution` for phase 2, and `scoring_output` is a directory that will be created.

6) Jupyter-notebook:

We provide a sample Jupyter notebook to help the participants getting started. Keep the basic structure of the template but replace the yellow place holders with relevant graphics and explanations:

jupyter-notebook README.ipynb.

7) From the starting kit to the challenge bundle:

The challenge data and starting kit are the basis for the Codalab challenge. Here is your **check list**:

- ✓ **DATA:** The challenge **data** should have this structure:

DirectoryName/	DirectoryName_train.data
	DirectoryName_train.solution
	DirectoryName_valid.data
	DirectoryName_valid.solution
	DirectoryName_test.data
	DirectoryName_test.solution
	DirectoryName_feat.name
	DirectoryName_label.name
	DirectoryName_feat.type
	DirectoryName_public.info
	DirectoryName_private.info

Put is in the directory **FILES/** at the same level as **starting_kit/**

In the starting kit, use as sample_data a small sample taken out of the training data.

- ✓ **STARTING KIT:** Your **starting kit** directory should have this structure:

- logo.jpg
- README.md.
- README.ipynb
- sample_data/
- sample_code_submission/
- ingestion_program/
- scoring_program/
- html_files/

- ✓ **LOGO:** Prepare a **logo** (a small square jpg file), call it **logo.jpg**.
- ✓ **TASK:** The contents of **ingestion_program/** should remain **unchanged**.
- ✓ **METRIC:** Declare your metric in **scoring_program/metric.txt**, corresponding to a metric of scikit-learn. **libscores.py**, OR the code of your own metric found in **my_metric.py**. Leave the rest of **scoring_program/** **unchanged**.
- ✓ **SAMPLE SUBMISSION:** You should have prepared a sample code submission in: **sample_code_submission/**. If it works, we will create automatically for you a **sample_result_submission/**. **Warning:** when you compile the bundle, the contents of this directory will be over-written by the new results generated by your sample code.
- ✓ **HTML FILES:** You should also prepare documentation files in HTML format in **html_files/**: You can edit them later so do not worry too much about them, you can keep the template.
- ✓ **README FILES:** README.ipynb (update with your own code) and README.md (keep its contents as is, except for title and description).

Compiling your challenge bundle

In the directory `starting_kit/utilities/`, you will find two files:

`competition.yaml`

`make_bundle.py`

This will allow you to compile your challenge bundle at the prompt.

Usage (1) with sample data only:

```
python make_bundle.py
```

Usage (2) with big data:

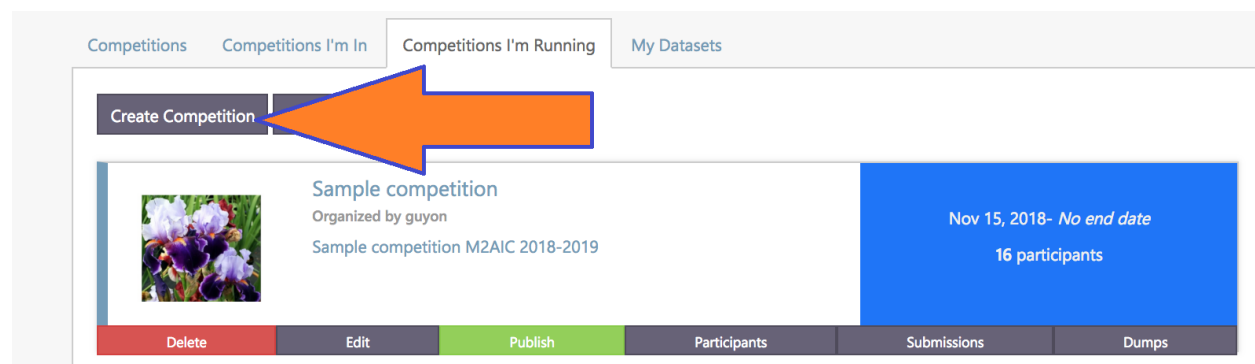
```
python make_bundle.py ../ ../.../FILES/iris
```

This creates a zip file `DataSet_bundle_date.zip` in the same directory.

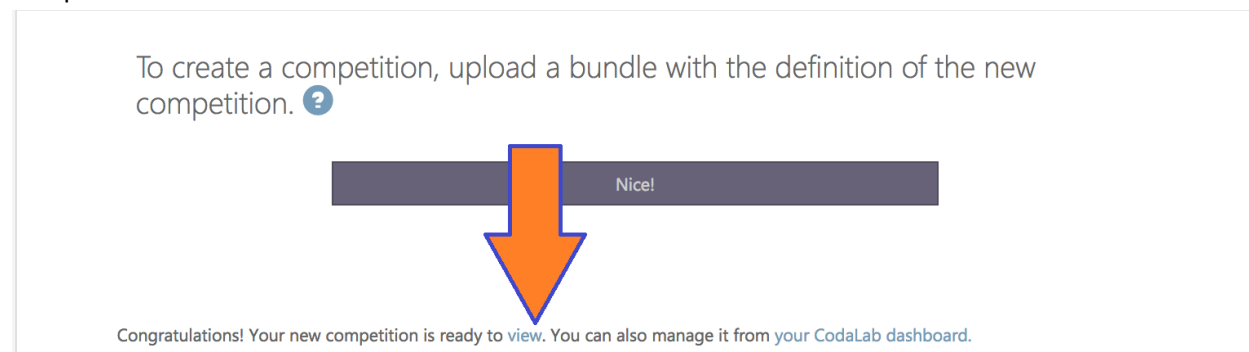
Upload the challenge bundle to Codalab

Upload `DataSet_bundle_date.zip` to <https://codalab.lri.fr/>.

To do this, use your account groupname@chalearn.org. Login and then go to “My Competitions” > “Competitions I’m Running” and click “Create Competition”.

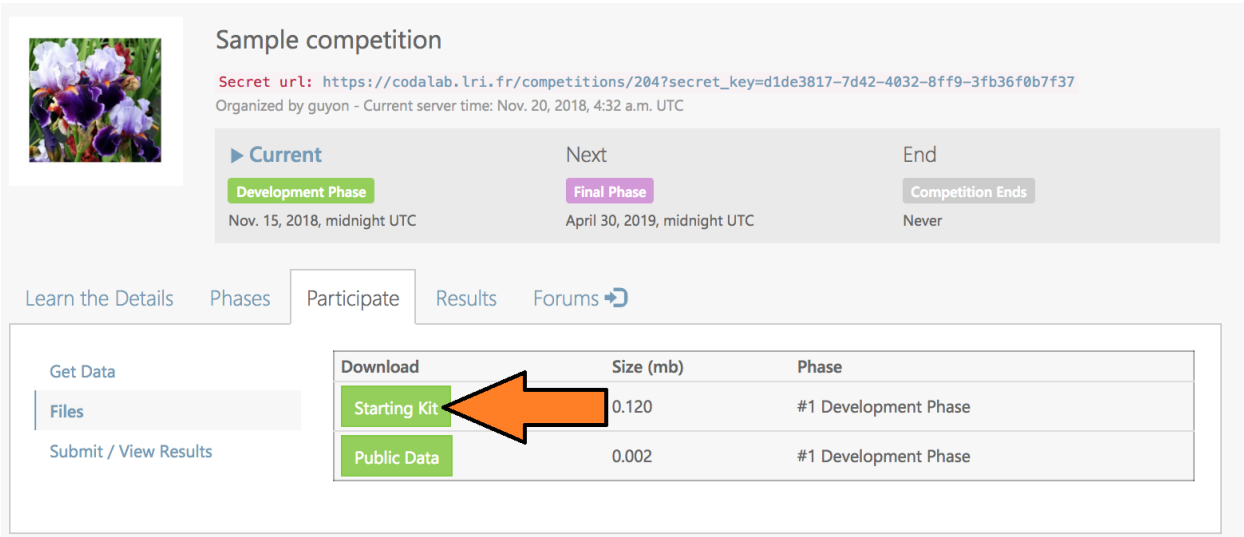


Once you upload your bundle, if all goes well, you can click on a microscopic link to view your competition:



If you upload the sample competition, this is what you get: <https://codalab.lri.fr/competitions/204>.

Do not upload a submission from your computer yet. First test the starting kit that is on the challenge website. Download it from the “Files” tab under “Participate”:



The screenshot shows the CodaLab competition interface for 'Sample competition'. It includes a competition image, a secret URL, and phase information. The 'Participate' tab is active, showing a table with download links for the 'Starting Kit' and 'Public Data'. An orange arrow points to the 'Starting Kit' link.

Sample competition

Secret url: https://codalab.lri.fr/competitions/204?secret_key=d1de3817-7d42-4032-8ff9-3fb36f0b7f37
Organized by guyon - Current server time: Nov. 20, 2018, 4:32 a.m. UTC

Current	Next	End
Development Phase Nov. 15, 2018, midnight UTC	Final Phase April 30, 2019, midnight UTC	Competition Ends Never

Learn the Details Phases **Participate** Results Forums ➔

Get Data

Files

Submit / View Results

Download	Size (mb)	Phase
Starting Kit	0.120	#1 Development Phase
Public Data	0.002	#1 Development Phase

In the same place, you can also download your “big” data (the public part only, not the validation and test labels).

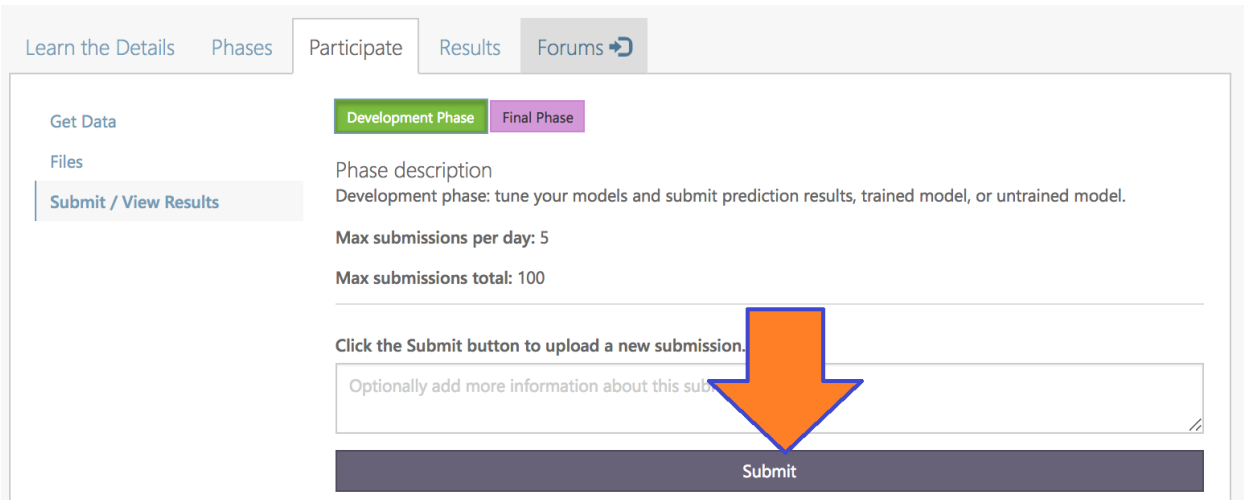
Then in your starting kit, locate your three sample submissions:

sample_code_submission.zip

sample_result_submission.zip

sample_trained_submission.zip

Upload them one by one to CodaLab to check that they work:



The screenshot shows the 'Participate' tab with the 'Submit / View Results' button highlighted. Below the phase information, there is a description of the development phase, submission limits, and a 'Submit' button. An orange arrow points to the 'Submit' button.

Learn the Details Phases **Participate** Results Forums ➔

Get Data

Files

Submit / View Results

Development Phase **Final Phase**

Phase description
Development phase: tune your models and submit prediction results, trained model, or untrained model.

Max submissions per day: 5
Max submissions total: 100

Click the Submit button to upload a new submission.

Optionally add more information about this submission.

Submit

Then look at the results on the leaderboard:

Learn the Details Phases Participate **Results** Forums ➔

Development Phase
Final Phase

Phase description
Development phase: tune your models and submit prediction results, trained model, or untrained model.

Max submissions per day: 5
Max submissions total: 100

Download CSV
Download all submissions on leaderboard

RESULTS						
#	User	Entries	Date of Last Entry	Prediction score ▲	Duration ▲	Detailed Results
1	guyon	4	11/20/18	0.7333 (1)	0.00 (1)	View

To do this you will need to go to **Participate>View/Submit results**.

WARNING: you will get errors if the correct docker is not specified. See page 11.

Large datasets

If you have a very large dataset you can compile your bundle with sample data only and the replace the baby datasets by the real ones.

1) Separate the target values from the rest and create 3 archives:

- **DirectoryName_input_data.zip**, containing:

DirectoryName_train.data
DirectoryName_train.solution
DirectoryName_valid.data
DirectoryName_test.data
DirectoryName_feat.name
DirectoryName_label.name
DirectoryName_feat.type
DirectoryName_public.info

- **DirectoryName_reference_validation_data.zip**, containing:

DirectoryName_valid.solution
DirectoryName_public.info
DirectoryName_private.info

- **DirectoryName_reference_test_data.zip**, containing:

DirectoryName_test.solution
DirectoryName_public.info
DirectoryName_private.info

go to **My Competitions>My Datasets>Create Dataset**. Fill out the form:

Dataset Create

Create Dataset

Name:
DataName_input_data

Description:
My large dataset

Type:
Input Data

Single new dataset Many existing datasets bundled into one new dataset

Data file:
Choose File No file chosen

Upload

Upload the file **DataName_input_data.zip**. After you upload, you should see your dataset in the data table. The KEY can be used to refer to it from the YAML file, both for public_data and input_data.

input_data: dac49905-dda0-4857-922a-02ca957ec8fd

public_data: dac49905-dda0-4857-922a-02ca957ec8fd

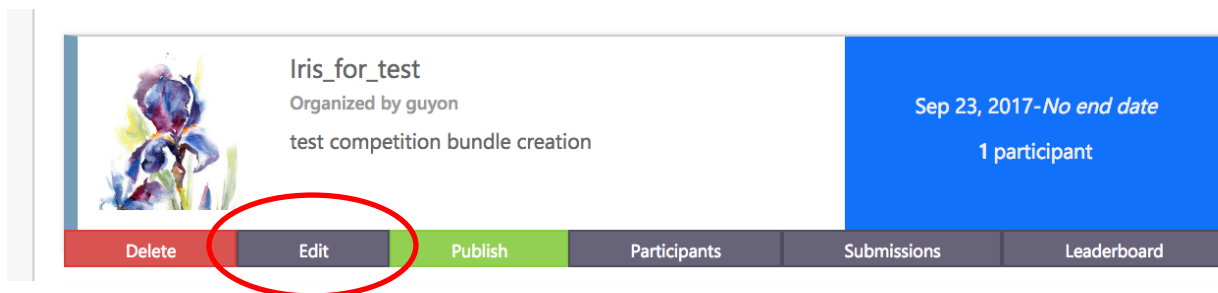
NAME	KEY	DESCRIPTION	TYPE	ACTIONS	MULTI-DELETE? (UNSAFE!)
DataName_input_data	dac49905-dda0-4857-922a-02ca957ec8fd	None	Input Data	Download Delete	<input type="checkbox"/>

Do the same thing for **DataName_reference_data.zip**. IMPORTANT, use the correct data Type (**input_data** or **reference_data**) in the form when you create the dataset. You can then use the key obtained in the YAML file:

reference_data: 0d8c9c0e-7609-4314-bf1c-3154cdee3462

You need only to specify public_data in phase 1. Use the same input_data and reference_data keys in both phases.

You can also use the editor. In Competitions I'm running, find your competition and click "Edit":



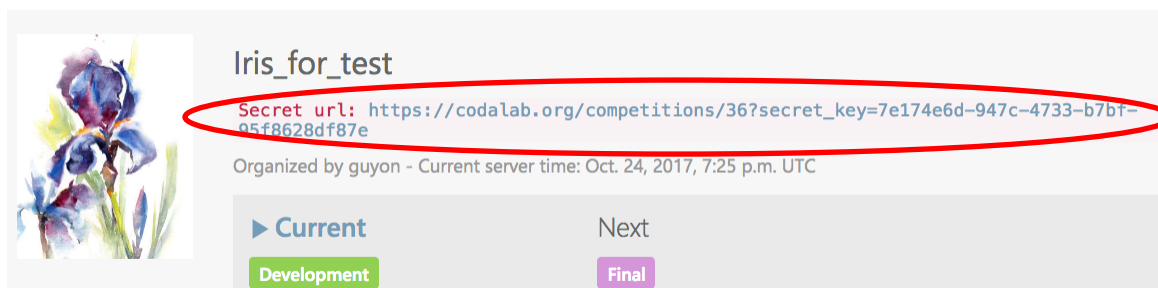
Find the menus for Input Data, Reference Data, and Public Data. Select the right dataset, as don't forget to **SAVE YOUR CHANGES**.

Input Data:

Reference Data:

Editing your challenge on Codalab

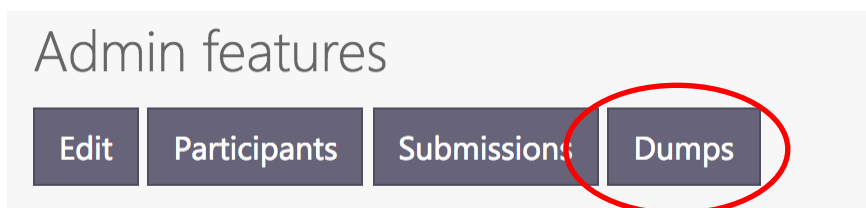
Once your challenge is uploaded to Codalab, your whole team can contribute to improving it. Before you make your challenge public (by clicking "Publish") you can share the secret URL with others.



But this will not allow them to edit the challenge. To give them this privilege, go to the editor and add their Codalab ID to the list of admins:

Admins:

The Codalab editor can be used to edit your challenge. However, try to minimize to avoid winding up with a broken site. Save your intermediate challenge bundles by using the **Dump** button.



Editing your challenge allows you to update your starting kit. You can upload it in my datasets, then point to it from the editor:

Starting Kit:

iris_baseline_1_72 uploaded by blue-chalearn.org

If your score is an “accuracy” (better is higher), you want the leaderboard sorted in descending order, but if it is an “error” (better is lower), you want it in ascending order. You can change that:

Sorting:

Descending

We want to run the code into Python 3 dockers. Make sure the correct docker is indicated:

Competition docker image:

codalab/codalab-legacy:py3

Using dockers

Codalab allows you to specify a docker in which you code will be running. Thus it is convenient to test your code in this docker ahead of time.

When you submit code to the Codalab platform, your code is executed inside a docker container. This environment can be exactly reproduced on your local machine by downloading the corresponding docker image. The `codalab/codalab-legacy:py3` docker, which can be pulled from <https://hub.docker.com/r/codalab/codalab-legacy/tags/> contains a large number of pre-loaded programs, including Python 3 and many libraries including scikit-learn. See <https://github.com/codalab/codalab-dockers/tree/master/legacy-py3> for details.

To run the starting kit inside the Codalab docker from a terminal:

Step 1) Create a temporary folder to put everything you will need (data and starting kit):

```
mkdir ~/aux
cp starting_kit.zip ~/aux
```

Copy any other data/code you want to that folder, as we will instruct you to map such local folder to be visible when you are within the docker. If you don't do such mapping procedure, you will not be able to see what you need from within the docker.

Step 2) Run the docker (and map such folder, described above, to `/home/aux`, inside the docker):

```
docker run -it -v ~/aux:/home/aux codalab/codalab-legacy:py3
```

The first time around it takes a while because docker must download the entire docker image. if everything goes well, you will be inside the docker and the prompt will be like below:

```
root@c74a8b1ccaf7:/#
```

Step 3) Run the starting kit. Inside the docker, go to “/home/aux” and unzip and run the starting kit.

```
# cd /home/aux
# unzip starting_kit.zip
# cd starting_kit
# python3 ingestion_program/ingestion.py sample_data sample_result_submission ingestion_program
sample_code_submission
# python3 scoring_program/score.py sample_data sample_result_submission scoring_output
# exit
```

WARNING: inside the docker, python 3.6 is called python3.

Step 4) Run the jupyter notebook.

Same thing as before, but add : **-p 8888:8888**

```
docker run -it -p 8888:8888 -v ~/aux:/home/aux codalab/codalab-legacy:py3
```

Go to a web browser and check that the notebook is running at <http://localhost:8888/>.

Then open README.pynb. WARNING: the default notebook kernel is Python 2, you’ll have to switch to Python 3. If you do not have a Python 3 kernel... you are in trouble. [Fixing it is complicated](#).

Modifying your docker

The docker we recommend contains a lot of libraries, but you may need more for your specific competition. You can then make changes to it [following these instructions](#). But AVOID THIS if you can because you will then have another docker as everyone else, this may confused the L2 students. Also, you will need to test that it works fine on Codalab. If you make a new docker, it must be posted on [Docker Hub](#). Preferable use a docker file to specify your changes and put them under revision control on Github.

How to use Git and Github

There are several good tutorials and cheat sheets on the Internet. You should start by cloning the repo of the starting kit, this is better than downloading the zip:

```
git clone https://github.com/madclam/m2aic2019.git
```

If you want to create another repo, [follow these instructions](#).