# RAMP
## DATA CHALLENGES WITH
## MODULARIZATION AND CODE SUBMISSION

# BALÁZS KÉGL

## Université Paris-Saclay / CNRS

université
PARIS-SACLAY

Paris-Saclay
Center for Data Science

# WHO AM I?
## Balázs Kégl

- Senior researcher **CNRS**

  - machine learning (20 years)
    interfacing with particle physics (10 years)

- Head of the **Paris-Saclay Center for Data Science**

  - interfacing with biology, economy, climatology, chemistry, etc. (4 years)

# Paris-Saclay Center for Data Science

**A multi-disciplinary initiative, building interfaces, matching people, helping them launching projects**

**345 affiliated researchers, 50 laboratories**

**Biology & bioinformatics**
IBISC/UEvry
LRI/UPSud
Hepatinov
CESP/UPSud-UVSQ-Inserm
IGM-I2BC/UPSud
MIA/Agro
MIAj-MIG/INRA
LMAS/Centrale

**Chemistry**
EA4041/UPSud

**Earth sciences**
LATMOS/UVSQ
GEOPS/UPSud
IPSL/UVSQ
LSCE/UVSQ
LMD/Polytechnique

**Economy**
LM/ENSAE
RITM/UPSud
LFA/ENSAE

**Neuroscience**
UNICOG/Inserm
U1000/Inserm
NeuroSpin/CEA

**Particle physics astrophysics & cosmology**
LPP/Polytechnique
DMPH/ONERA
CosmoStat/CEA
IAS/UPSud
AIM/CEA
LAL/UPSud

**Machine learning**
LRI/UPSud
LTCI/Telecom
CMLA/Cachan
LS/ENSAE
LIX/Polytechnique
MIA/Agro
CMA/Polytechnique
LSS/Supélec
CVN/Centrale
LMAS/Centrale
DTIM/ONERA
IBISC/UEvry
LIST/CEA

**Visualization**
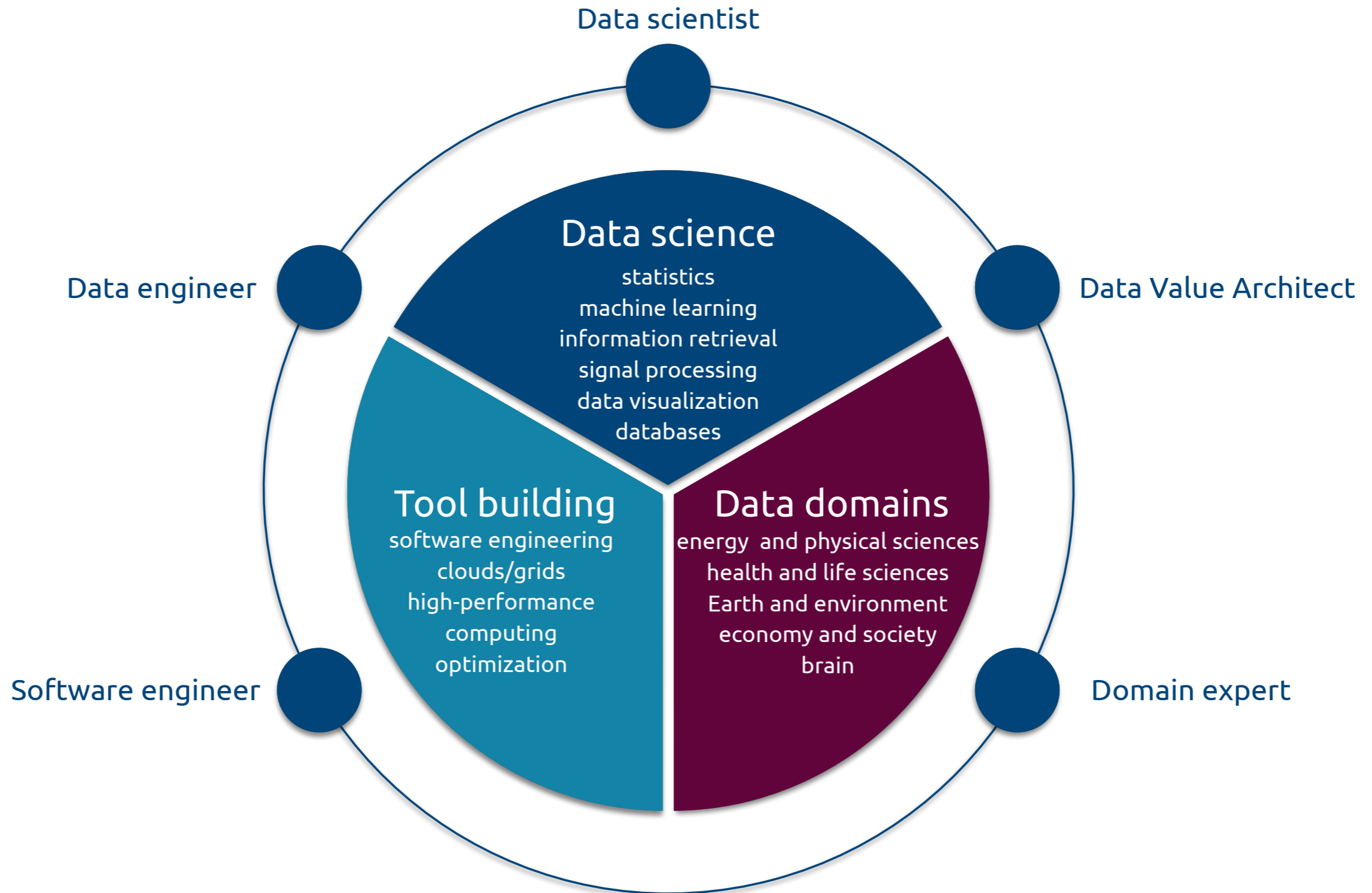INRIA
LIMSI

**Signal processing**
LTCI/Telecom
CMA/Polytechnique
CVN/Centrale
LSS/Supélec
CMLA/Cachan
LIMSI
DTIM/ONERA

**Statistics**
LMO/UPSud
LS/ENSAE
LSS/Supélec
CMA/Polytechnique
LMAS/Centrale
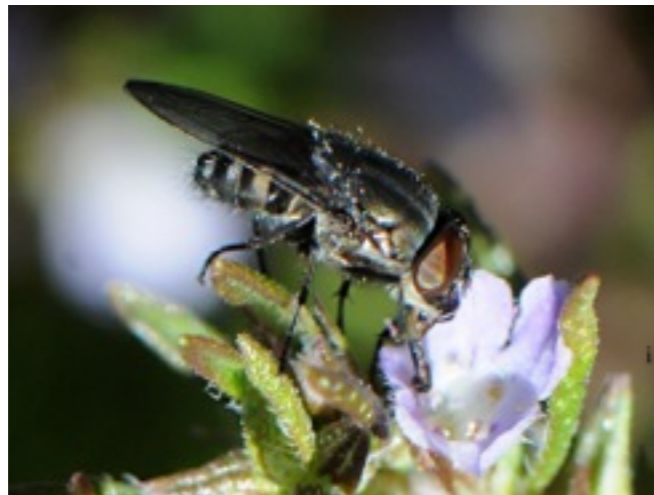MIA/AgroParisTech

B. Kégl (CNRS)

3

# The data science ecosystem

## https://medium.com/@balazskegl

# WHAT DOES MACHINE LEARNING DO

- **Classification** problem *y = f(x)*

*x*



*f* → *y* 'Stomorhina'

*x*



*f* → *y* 'Scaeva'

# WHAT DOES MACHINE LEARNING DO

x

y

**FREEDOMJUNKSHUN.COM**

Says "FBI and Texas State Police are both confirming that Raymond Peter Littleberry, the man accused of shooting up a Texas church, was an avid atheist on the payroll of the Democrat National Committee."

— *PolitiFact National* on Monday, November 6th, 2017

PANTS ON FIRE!
POLITIFACT
TRUTH-O-METER™

Not on DNC payroll

**DAVID PERDUE**

The diversity visa lottery is "plagued by fraud."

— *PolitiFact National* on Monday, November 6th, 2017

MOSTLY TRUE
POLITIFACT
TRUTH-O-METER™

Some fraud rings extract money from lottery winners

**BARACK OBAMA**

"Eight in 10 people this year can find plans for $75 a month or less."

— *PolitiFact National* on Monday, November 6th, 2017

HALF TRUE
POLITIFACT
TRUTH-O-METER™

Under 4 percent of Americans

# ML IN NATURAL SCIENCES

- Inference: to invert the generative model

  - "predict" a particle, detect an anomaly, infer a parameter y, from observation x

- Generation: to replace expensive simulations

  - "learn" a GEANT4 simulation with a neural net

- Hypothesis generation: to replace theoreticians :)

  - learn, represent structural knowledge and generate novelty in model space, e.g., molecule generation in drug discovery

# RAMP is a **tool** for

1. **Collaborative prototyping**
2. **Teaching aid**
3. **Data science process management**

# Funded by Université Paris-Saclay and CNRS

## Team



Balázs Kégl

Alex Gramfort

Akin Kazakçi

Mehdi Cherti

Yohann Sitruk

Guillaume Lemaître

Alexandre Boucaud

Joris Van den Bossche

## Alumni

Djalel Benbouzid

Camille Marini

# RAMP.STUDIO
## DATA CHALLENGE WITH CODE SUBMISSION

# Code submission

1. lets us deliver a **working prototype**
2. lets the participants **collaborate**
3. makes the **backend challenging** to run (cloud management)

# RAPID ANALYTICS AND MODEL PROTOTYPING



**functional data**, **time series**, **data augmentation**, **deep learning**, **learning on simulations**, **nonstandard and multi-objective losses**

université PARIS-SACLAY

Paris-Saclay
Center for Data Science

# CURRENT RAMPS

RAMP

Hi Balazs! ▾

- **Pollenating insect classification (209 classes)**
  - La Paillasse / Futur en Seine, number of participants = **28**, number of submissions = **13**, combined score = **0.831**, click here for score vs time plot
- **Titanic survival classification**
  - DSSP 6 2016/17 2, number of participants = **31**, number of submissions = **34**, combined score = **0.87**, click here for score vs time plot
  - Entry exam to deep learning tutorial, number of participants = **35**, number of submissions = **21**, combined score = **0.86**, click here for score vs time plot
  - Ecole des Mines 2016/17, number of participants = **125**, number of submissions = **144**, combined score = **0.89**, click here for score vs time plot
- **Pollenating insect classification (18 classes)**
  - Polytechnique MAP583/MAP542 2016/17, number of participants = **166**, number of submissions = **114**, combined score = **0.959**, click here for score vs time plot
  - DSSP5 2017, number of participants = **15**, number of submissions = **24**, combined score = **0.93**, click here for score vs time plot
- **Particle tracking in the LHC ATLAS detector**
  - initial single-day RAMP 2017, number of participants = **55**, number of submissions = **60**, combined score = **0.97**, click here for score vs time plot
- **El Nino forecast**
  - single-day RAMP at Climate Informatics Workshop 2015; Saclay Data Camp 2016/17, number of participants = **160**, number of submissions = **138**, combined score = **0.389**, click here for score vs time plot
- **Arctic sea ice forecast**
  - single-day RAMP at Climate Informatics Workshop 2016, number of participants = **46**, number of submissions = **83**, combined score = **0.31**, click here for score vs time plot
  - Polytechnique MAP542 2016/17, number of participants = **20**, number of submissions = **52**, combined score = **0.268**, click here for score vs time plot
  - Polytechnique MAP583 2016/17, number of participants = **123**, number of submissions = **252**, combined score = **0.259**, click here for score vs time plot
- **Number of air passengers prediction**
  - DSSP4/5 2016, number of participants = **95**, number of submissions = **242**, combined score = **0.236**, click here for score vs time plot
  - DSSP6 2017, number of participants = **23**, number of submissions = **59**, combined score = **0.268**, click here for score vs time plot
- **Drug classification and concentration estimation from Raman spectra**
  - Polytechnique MAP583 2016/17, number of participants = **125**, number of submissions = **258**, combined score = **0.048**, click here for score vs time plot
  - initial single-day RAMP 2016; Saclay Data Camp 2016/17, number of participants = **242**, number of submissions = **554**, combined score = **0.027**, click here for score vs time plot
  - Ecole des Mines 2016/17, number of participants = **124**, number of submissions = **560**, combined score = **0.023**, click here for score vs time plot
- **Detecting anomalies in the LHC ATLAS detector**
  - Polytechnique MAP542 2016/17, number of participants = **29**, number of submissions = **47**, combined score = **0.865**, click here for score vs time plot
  - Polytechnique MAP583 2016/17, number of participants = **133**, number of submissions = **275**, combined score = **0.899**, click here for score vs time plot
  - initial single-day RAMP 2016, number of participants = **49**, number of submissions = **19**, combined score = **0.677**, click here for score vs time plot
- **Epidemium cancer mortality rate prediction (2nd RAMP)**
  - initial single-day RAMP 2016, number of participants = **39**, number of submissions = **46**, combined score = **21.79**, click here for score vs time plot
  - Polytechnique MAP583 2016/17, number of participants = **128**, number of submissions = **192**, combined score = **18.59**, click here for score vs time plot
  - Polytechnique MAP542 2016/17, number of participants = **22**, number of submissions = **57**, combined score = **19.31**, click here for score vs time plot

université PARIS-SACLAY

Paris-Saclay
**Center for Data Science**

# DATA SCIENCE THEMES

## Data science themes

- **classification**
  - Iris classification
  - Detecting anomalies in the LHC ATLAS detector
  - Drug classification and concentration estimation from Raman spectra
  - Titanic survival classification
  - Pollenating insect classification (18 classes)
  - Pollenating insect classification (209 classes)
- **convolutional networks**
  - Pollenating insect classification (18 classes)
  - Pollenating insect classification (209 classes)
- **external data**
  - Number of air passengers prediction
- **feature engineering**
  - El Nino forecast
  - Arctic sea ice forecast
  - Drug classification and concentration estimation from Raman spectra
  - Detecting anomalies in the LHC ATLAS detector
- **forests**
  - Iris classification
  - Detecting anomalies in the LHC ATLAS detector
  - Titanic survival classification
  - Boston housing price regression
  - El Nino forecast
  - Arctic sea ice forecast
  - Number of air passengers prediction
  - Epidemium cancer mortality rate prediction (2nd RAMP)
- **functional data**
  - Drug classification and concentration estimation from Raman spectra
- **image data**
  - Pollenating insect classification (18 classes)
  - Pollenating insect classification (209 classes)
  - El Nino forecast

- **missing data**
  - Epidemium cancer mortality rate prediction (2nd RAMP)
  - Titanic survival classification
- **neural networks (deep learning)**
  - Drug classification and concentration estimation from Raman spectra
  - Pollenating insect classification (18 classes)
  - Pollenating insect classification (209 classes)
- **regression**
  - Boston housing price regression
  - El Nino forecast
  - Arctic sea ice forecast
  - Number of air passengers prediction
  - Drug classification and concentration estimation from Raman spectra
  - Epidemium cancer mortality rate prediction (2nd RAMP)
- **small data**
  - Drug classification and concentration estimation from Raman spectra
  - Epidemium cancer mortality rate prediction (2nd RAMP)
  - Detecting anomalies in the LHC ATLAS detector
  - El Nino forecast
  - Arctic sea ice forecast
  - Number of air passengers prediction
  - Particle tracking in the LHC ATLAS detector
- **supervised clustering (unsupervised classification)**
  - Particle tracking in the LHC ATLAS detector
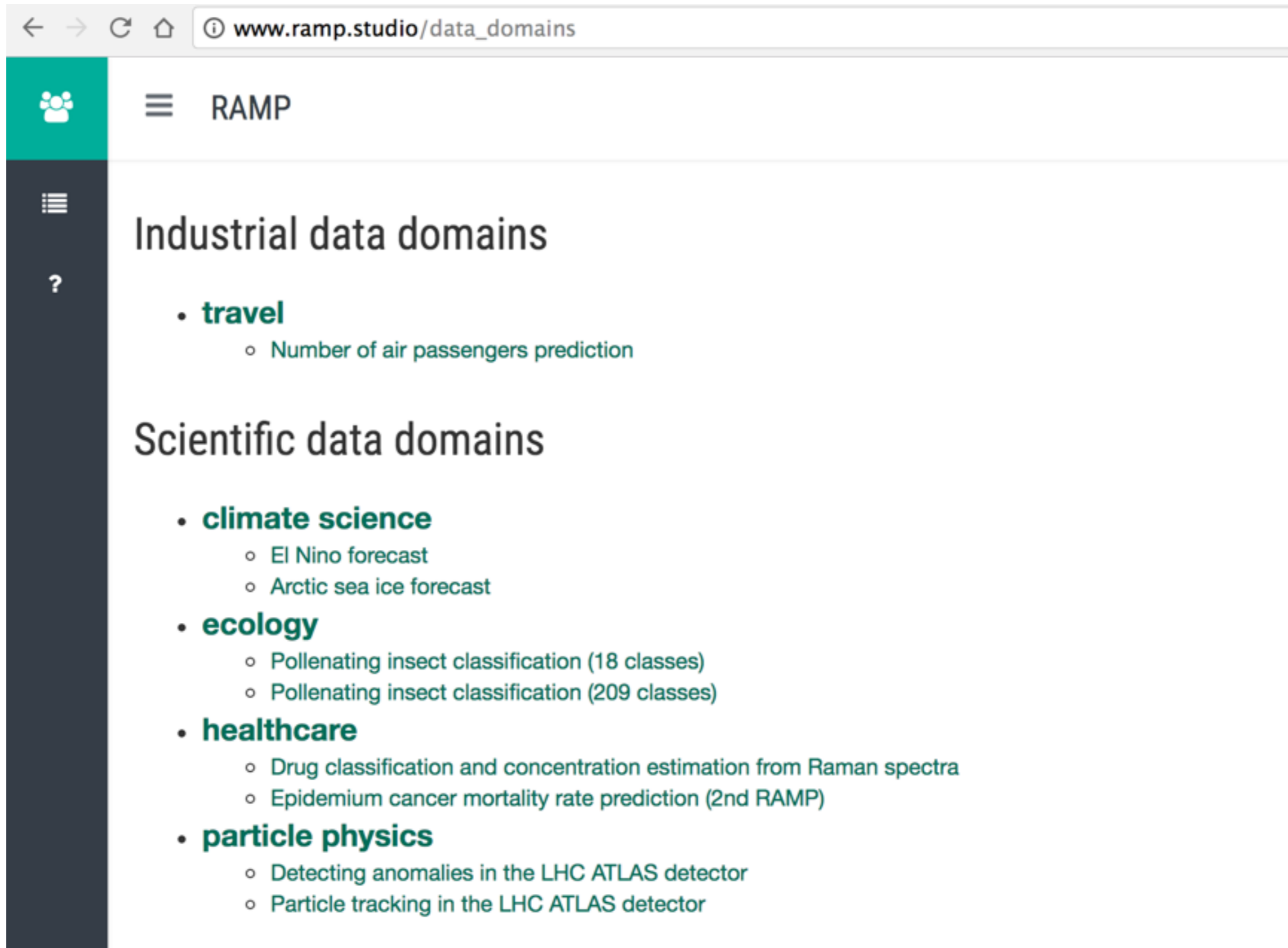- **tabular data**
  - Iris classification
  - Detecting anomalies in the LHC ATLAS detector
  - Titanic survival classification
  - Boston housing price regression
  - Number of air passengers prediction
  - Epidemium cancer mortality rate prediction (2nd RAMP)
- **time series forecasting**
  - El Nino forecast
  - Arctic sea ice forecast

# DATA DOMAINS

B. Kégl (CNRS)

# RAMP.STUDIO

## *DATA CHALLENGE WITH CODE SUBMISSION*

**20+ challenges**

**40+ events**

**1200+ users**

**7000+ predictive models**

# RAMP.STUDIO
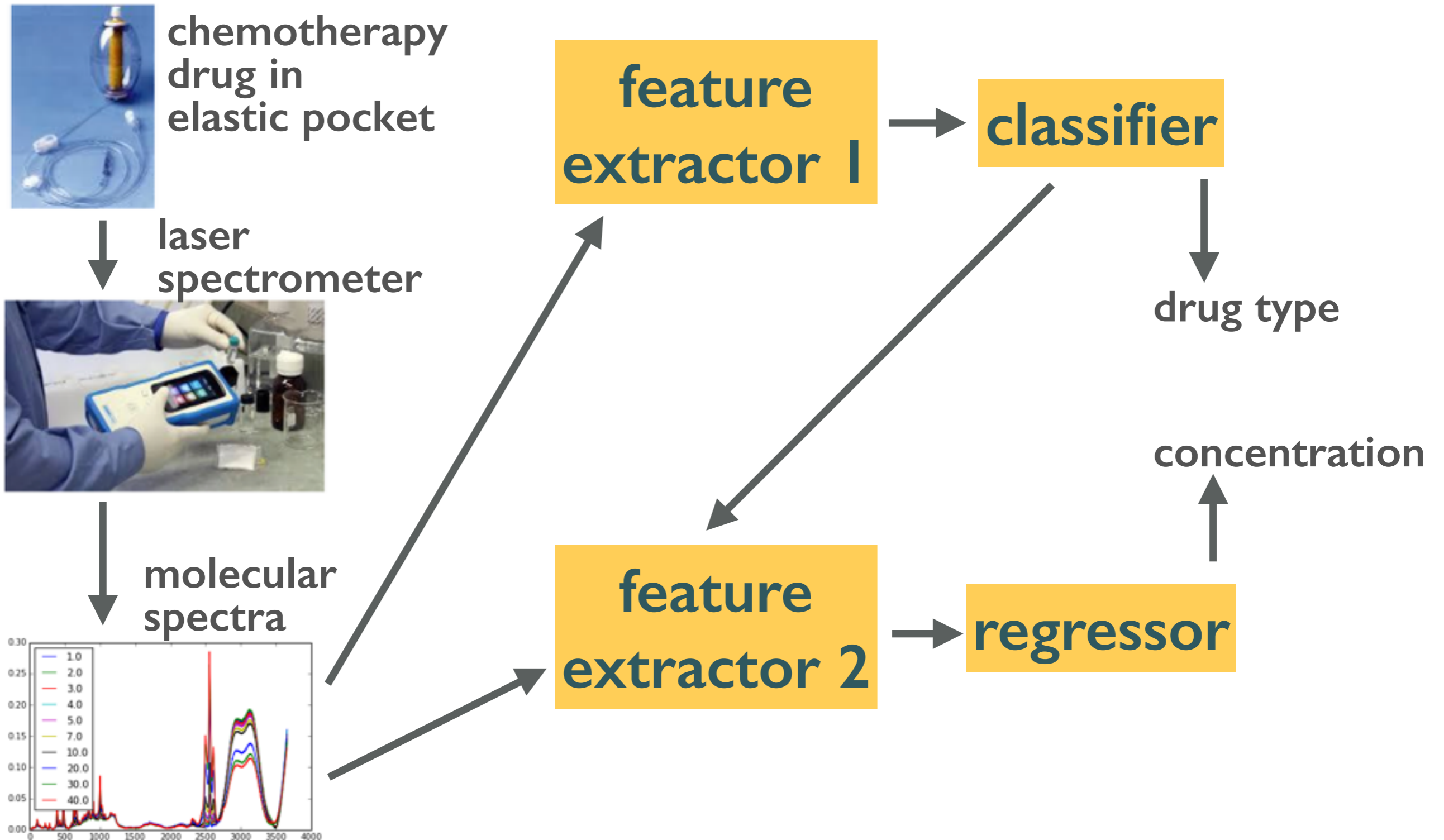## *DATA CHALLENGE WITH CODE SUBMISSION*

**12 hackathons**
**6 remote data challenges**
**11 course data camps**

# CLASSIFYING AND REGRESSING ON MOLECULAR SPECTRA



chemotherapy drug in elastic pocket

laser spectrometer

molecular spectra

feature extractor 1 → classifier → drug type

classifier → feature extractor 2 → regressor → concentration

université PARIS-SACLAY

Paris-Saclay Center for Data Science

# Classifying and quantifying monoclonal antibody preparations for cancer therapy using machine learning

Laetitia Le [ab], Camille Marini [ce], Alexandre Gramfort [cfg],
David Nguyen [a], Mehdi Cherti [ch], Sana Tfaili [b], Ali
Tfayli [b], Arlette Baillet-Guffroy [b], Eric Caudron [ab], Balázs
Kégl [ch]

[a] European Georges Pompidou Hospital (AP-HP), Pharmacy
department, Paris, France
[b] Lip(Sys) Chimi Analytique Pharmaceutique, Univ. Paris-Sud,
Universit Paris Saclay, F92290 Chatenay-Malabry, France
(EA4041 Groupe de Chimie Analytique de Paris Sud)
[c] Center of Data Science, Université Paris-Saclay
[d] Université Paris-Sud
[e] CMAP, Ecole Polytechnique, Palaiseau, France
[f] INRIA, Parietal team, Saclay, France
[g] LTCI, Télécom ParisTech
[h] LAL, CNRS, France

université
PARIS-SACLAY

Paris-Saclay
Center for Data Science

# FORECASTING EL NINO SIX MONTHS AHEAD

université PARIS-SACLAY

Paris-Saclay
Center for Data Science

# The power of the (collaborating) crowd



Hep detector anomalies test scores

# THE DYNAMICS OF COLLABORATION



Drug spectra submissions

starting kit

inventors

early influencers

the crowd

- ○ open phase
- ○ closed phase
- ● contributivity
- ● historical contributivity

# OPEN PHASE LETS PARTICIPANTS CATCH UP

## *THE GOAL OF TEACHING*



Hep detector anomalies test score histograms

# COMMUNICATION AND REUSE



Hep detector anomalies submissions

Legend:
- credited
- looked at
- closed phase
- historical contributivity
- contributivity
- open phase

team: etienne.boursier
subm: combine_features
score: 0.896
contr: 2
timestamp: 7d 7h 26m

submissions

# CLASSIFY POLLENATING INSECTS

**https://www.ramp.studio/events/pollenating_insects_3_JNI_2017**

**4.5K€** for the **competitive** phase
**3K€** for the **collaborative** phase
**50 GPU hours** per participant

# DETECTING MARS CRATERS

# UPCOMING CHALLENGE
# PREDICT AUTISM FROM BRAIN SCANS

# UPCOMING CHALLENGE(S)
# ASTROIMAGING



HSC data

Detection

↓

Deblending

↓                    ↓

Photometry          Morphology

# THE DATA FLOW

data connectors

$X$

predictive workflow

FE → CLF → CAL

$y_{pred}$

full automation production

dashboard decision support

# Before solving the problem, **set it up** (even put it into production)

# Setting up the RAMP
## ~~is~~ was
## long and hard.

# Separate workflow **building** and workflow **optimization**

# RAMP

## ~~Data-challenges-with~~

## ~~modularization~~ ~~and~~ ~~code-submission~~

## Frugal data science

## process management

# Balázs Kégl

## Université Paris-Saclay / CNRS

Balázs Kégl

Université Paris-Saclay / CNRS

universite
PARIS-SACLAY

Paris-Saclay
Center for Data Science

# RAMP-WORKFLOW & RAMP-KITS

- toolkit: https://github.com/paris-saclay-cds/ramp-workflow

  - for **designing workflows**

  - set of ready-made **metrics, workflows, CV schemes,** data readers

  - unique command-line **test script**

- examples: https://github.com/**ramp-kits**

  - a zoo of **problems**, **experiments**, **workflows**

  - (at least) one **initial solution**

# A SINGLE SCRIPT TO DEFINE THE BUNDLE

```python
27
28
29  def get_cv(X, y):
30      unique_replicates = np.unique(X['replicate'])
31      r = np.arange(len(X))
32      for replicate in unique_replicates:
33          train_is = r[(X['replicate'] != replicate).values]
34          test_is = r[(X['replicate'] == replicate).values]
35          yield train_is, test_is
36
37
38  def _read_data(path, f_name):
39      data = pd.read_csv(os.path.join(path, 'data', f_name))
40      y_array = data[_target_column_name]
41      X_df = data.drop([_target_column_name], axis=1)
42      return X_df, y_array
43
44
45  def get_train_data(path='.'):
46      f_name = 'train.csv.gz'
47      return _read_data(path, f_name)
48
49
50  def get_test_data(path='.'):
51      f_name = 'test.csv.gz'
52      return _read_data(path, f_name)
```
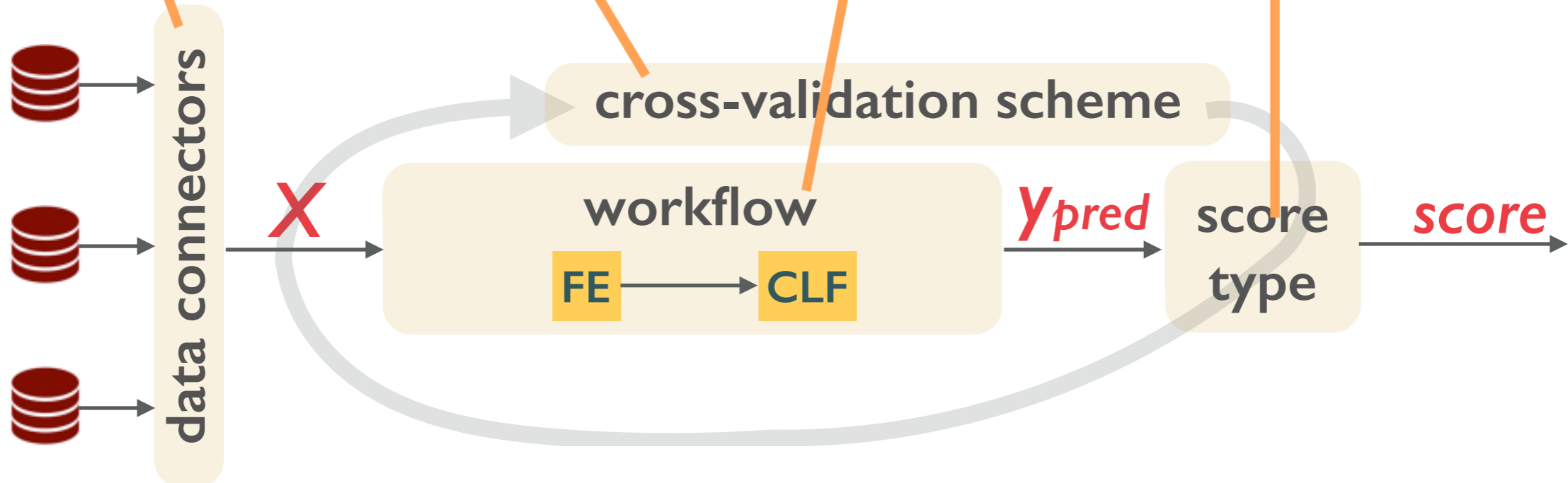
```python
1   import os
2   import numpy as np
3   import pandas as pd
4   import rampwf as rw
5
6   problem_title = \
7       'Cell population identification from single-cell mass cytometry data'
8   _target_column_name = 'cell type'
9   _prediction_label_names = [
10      'B-cell Frac A-C (pro-B cells)', 'Basophils', 'CD4 T cells', 'CD8 T cells',
11      'CLP', 'CMP', 'Classical Monocytes', 'Eosinophils', 'GMP', 'HSC',
12      'IgD- IgMpos B cells', 'IgDpos IgMpos B cells', 'IgM- IgD- B-cells',
13      'Intermediate Monocytes', 'MEP', 'MPP', 'Macrophages', 'NK cells',
14      'NKT cells', 'Non-Classical Monocytes', 'Plasma Cells', 'gd T cells',
15      'mDCs', 'pDCs']
16  # A type (class) which will be used to create wrapper objects for y_pred
17  Predictions = rw.prediction_types.make_multiclass(
18      label_names=_prediction_label_names)
19  # An object implementing the workflow
20  workflow = rw.workflows.FeatureExtractorClassifier()
21
22  score_types = [
23      rw.score_types.BalancedAccuracy(name='bac', precision=3),
24      rw.score_types.Accuracy(name='acc', precision=3),
25      rw.score_types.NegativeLogLikelihood(name='nll', precision=3),
26  ]
```

**data connectors** → $X$ → **cross-validation scheme**

**workflow** FE → CLF → $y_{pred}$ → **score type** → *score*

37

# A SINGLE EXECUTABLE TO TEST THE SUBMISSIONS

- Keep your different submissions in a **simple file structure**

- Communicate them on **git**

- Execute them also from the **notebook**

```
silver6:mouse_cytometry kegl$ ramp_test_submission
Testing Cell population identification from single-cell mass cytometry data
Reading train and test files from ./data ...
Reading cv ...
Training ./submissions/starting_kit ...
CV fold 0
        train bac = 0.042
        valid bac = 0.042
        test bac = 0.042
        train acc = 0.427
        valid acc = 0.416
        test acc = 0.396
        train nll = 1.715
        valid nll = 1.71
        test nll = 1.779
CV fold 1
        train bac = 0.042
        valid bac = 0.046
        test bac = 0.042
        train acc = 0.415
        valid acc = 0.453
        test acc = 0.396
        train nll = 1.729
        valid nll = 1.657
        test nll = 1.775
CV fold 2
        train bac = 0.042
        valid bac = 0.042
        test bac = 0.042
        train acc = 0.408
        valid acc = 0.471
        test acc = 0.394
        train nll = 1.738
        valid nll = 1.61
        test nll = 1.772
CV fold 3
        train bac = 0.042
        valid bac = 0.043
        test bac = 0.042
        train acc = 0.448
        valid acc = 0.357
        test acc = 0.396
        train nll = 1.655
        valid nll = 1.915
        test nll = 1.789
--------------------------
train bac = 0.042 ± 0.0001
train acc = 0.425 ± 0.0152
train nll = 1.709 ± 0.0325
valid bac = 0.043 ± 0.0016
valid acc = 0.424 ± 0.0437
valid nll = 1.723 ± 0.1165
test bac = 0.042 ± 0.0001
test acc = 0.395 ± 0.0006
test nll = 1.779 ± 0.0062
```

# You can

1. **Participate** in upcoming RAMPs
2. Use RAMP in **teaching** or **training**
3. Use the toolkit for **your own workflows**
4. Submit it to us if you want to **run a data challenge**

frontend:

**www.ramp.studio**

toolkit:

**github.com/paris-saclay-cds/ramp-workflow**

examples:

**github.com/ramp-kits**

slack:

**ramp-studio.slack.com**

blogs:

**medium.com/@balazskegl**

mail:

**balazs.kegl@gmail.com**