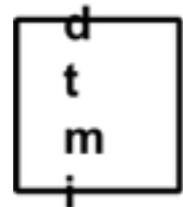


A data science observatory

Akin Kazakci, Mines ParisTech
Balazs Kégl, CNRS



design theory
and methods
for innovation



Team



Balázs Kégl
CNRS



Alexandre Gramfort
Télécom ParisTech



Akın Kazakçı
Mines ParisTech



Djalel Benbouzid
UPMC



Camille Marini
Télécom ParisTech



Mehdi Cherti
UP Saclay



Yohann Sitruk
Mines ParisTech

The research objective & questions

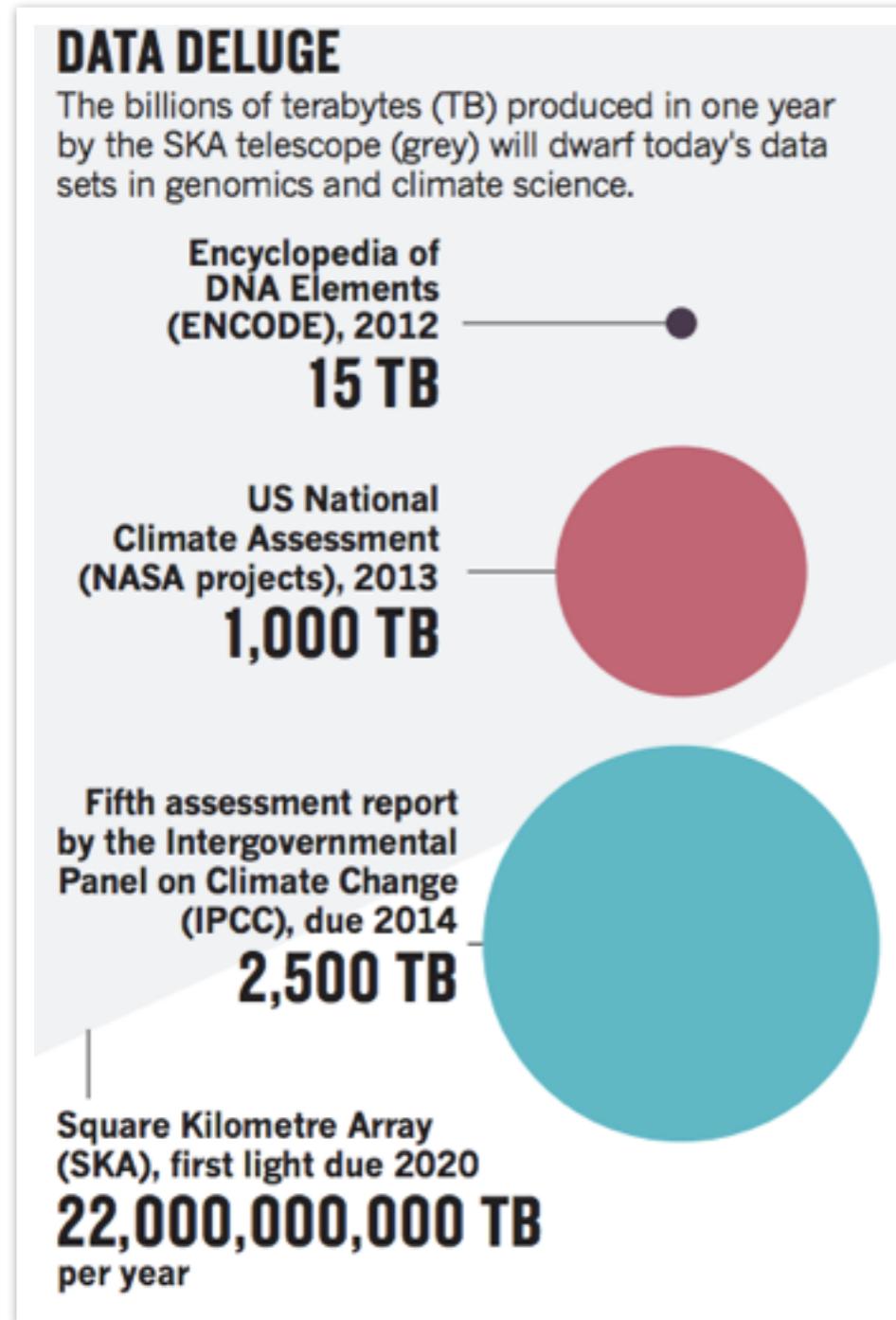
1

Enough with the chairs



- Design research is falling behind in dealing with **contemporary challenges**
 - Claim 1a: Too much in-breeding and repetition
 - Claim 1b: Huge amount of work is based on ideas from 80s'
- Design is **not about objects, but about reasoning**

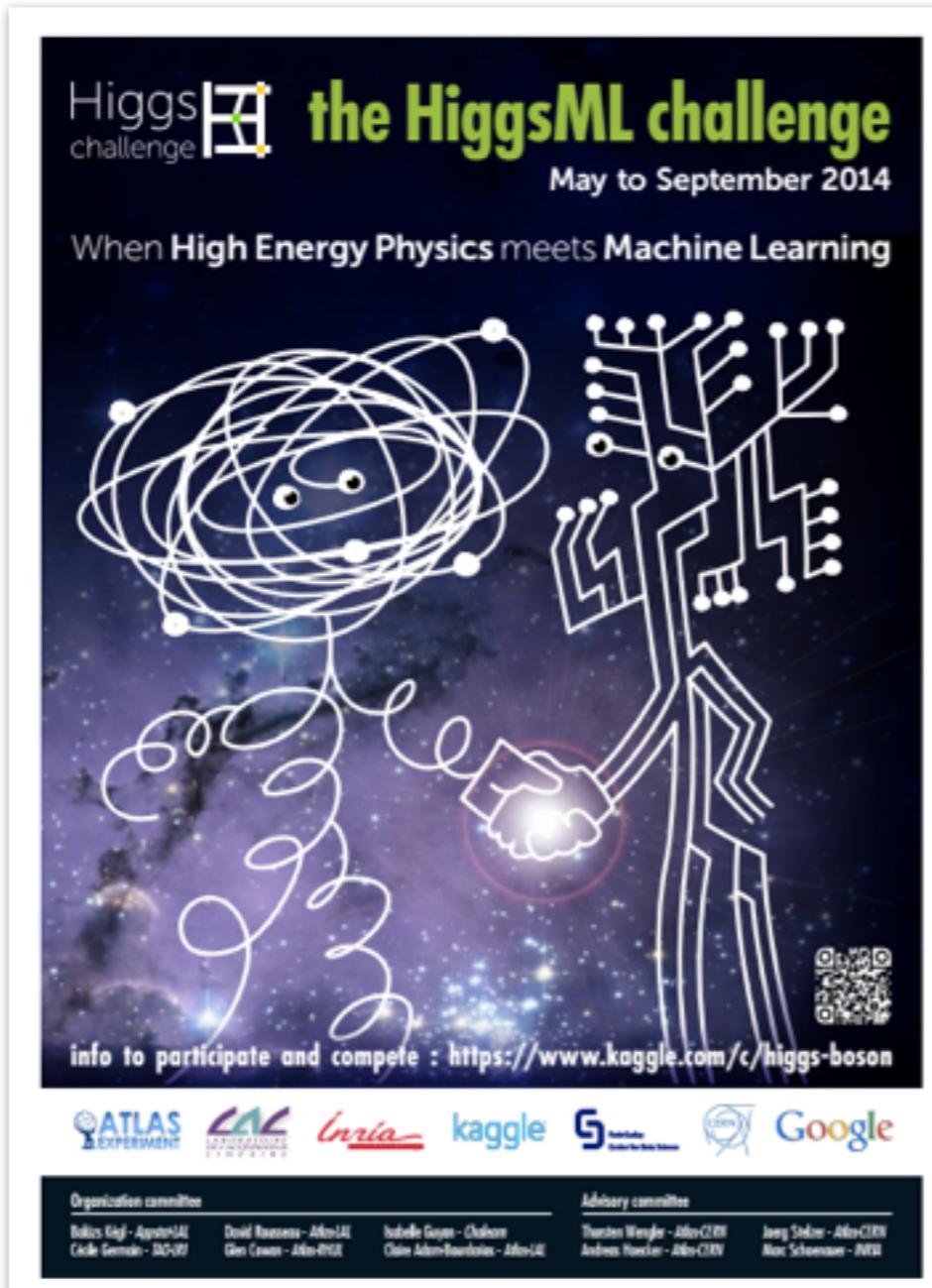
Revealing the potential of data: what role for design?



- Physics (Particle physics, Plasma physics, astrophysics...)
- Biology (Genetics, Epidemiology...)
- Chemistry
- Economics, Finance, Banking
- Manufacturing, Industrial Internet
- Internet of things, Connected Devices
- Social media
- Transport & Mobility
- ...

There is not enough data scientists to handle this much data

Last year: Crowdsourcing data challenges

A screenshot of the Kaggle website showing the results of the "Higgs Boson Machine Learning Challenge". The page header includes the Kaggle logo, navigation links for "Customer Solutions", "Competitions", "Community", and buttons for "Sign up" and "Login". The main title "Higgs challenge" is displayed, along with the text "Completed • \$13,000 1,785 teams". A large red circle highlights the number "1,785 teams". Below this, the text "Higgs Boson Machine Learning Challenge" and "Mon 12 May 2014 – Mon 15 Sep 2014 (21 days)" is shown. A large red "1785 teams" is overlaid on the page. A "Dashboard" button is visible, and the text "Private Leaderboard - Higgs Boson Machine Learning Challenge" is displayed. A note states "This competition has completed. This leaderboard reflects the final standings." and "See someone using multiple accounts? Let us know." A table shows the top three entries:

#	Δ1w	Team Name	model uploaded	*In the money	Score	Entries	Last Submission UTC (Best - Last Submission)
1	+4	Gábor Melis ± *			3.80581	110	Sun, 14 Sep 2014 09:10:04 (-0h)
2	-1	Tim Salimans ± *			3.78913	57	Mon, 15 Sep 2014 23:49:02 (-40.6d)
3	-	nhlxShaze ± *			3.78682	254	Mon, 15 Sep 2014 16:50:01 (-76.3d)

*Reasonable doubts about
the effectiveness of data
science contests*

Crowdsourcing /?/ Design

crowd

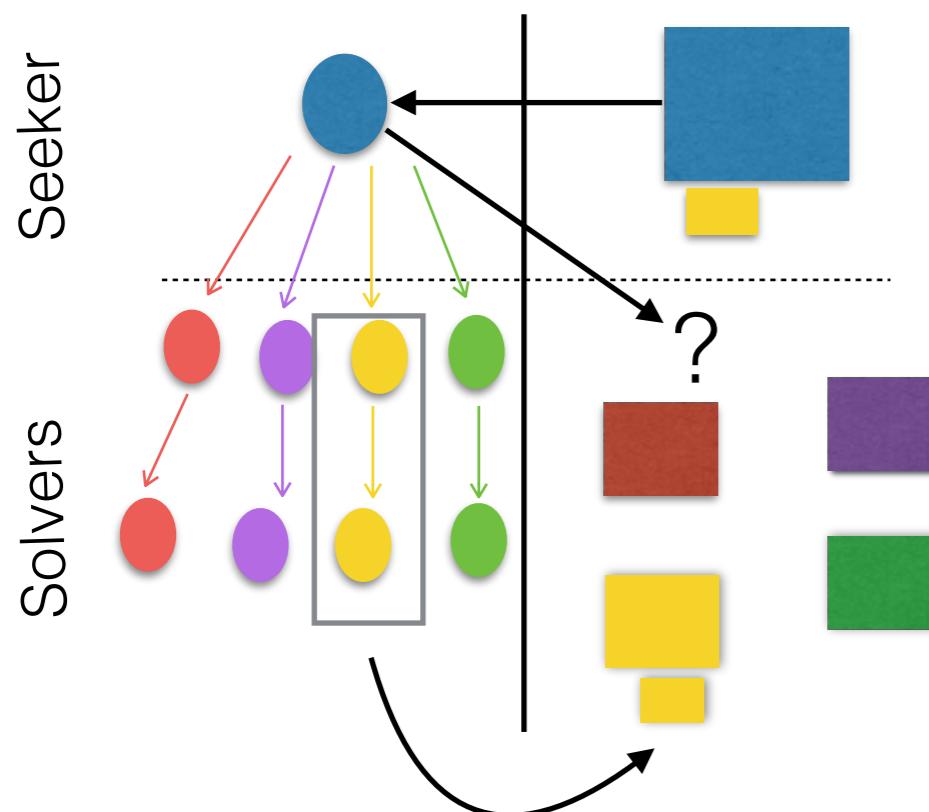
/kraʊd/ noun

1.a large number of people gathered together in a
disorganized or unruly way

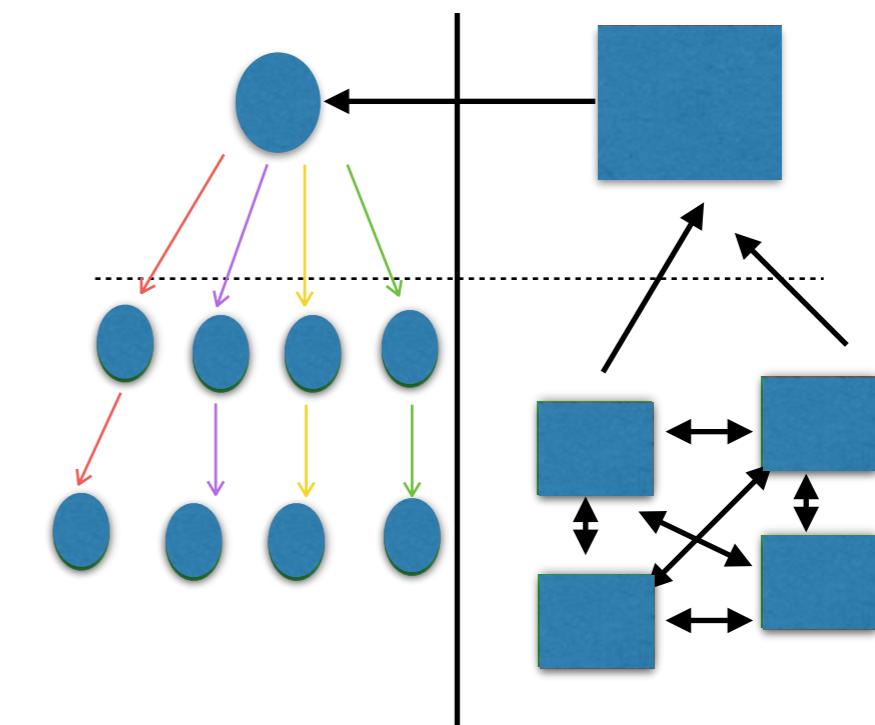
1. How to study the design process of a crowd?
2. How to manage the design process of a crowd?

Crowdsourcing: C-K dynamics

Crowdsourced contests

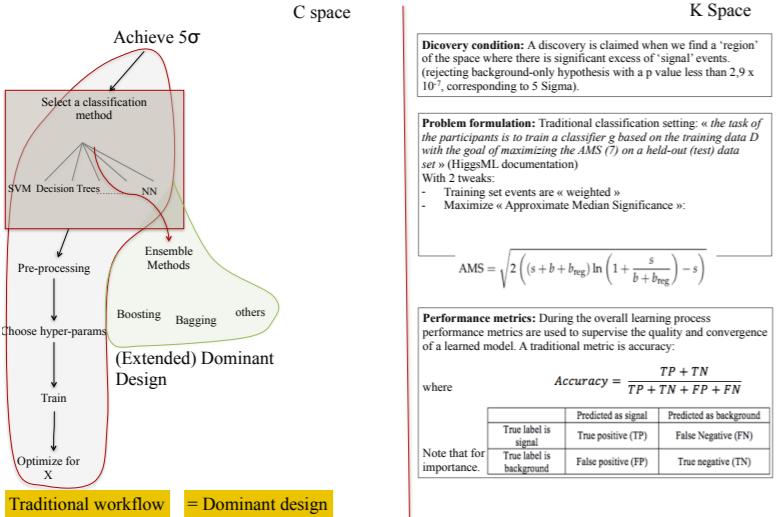


Crowdsourced collaboration

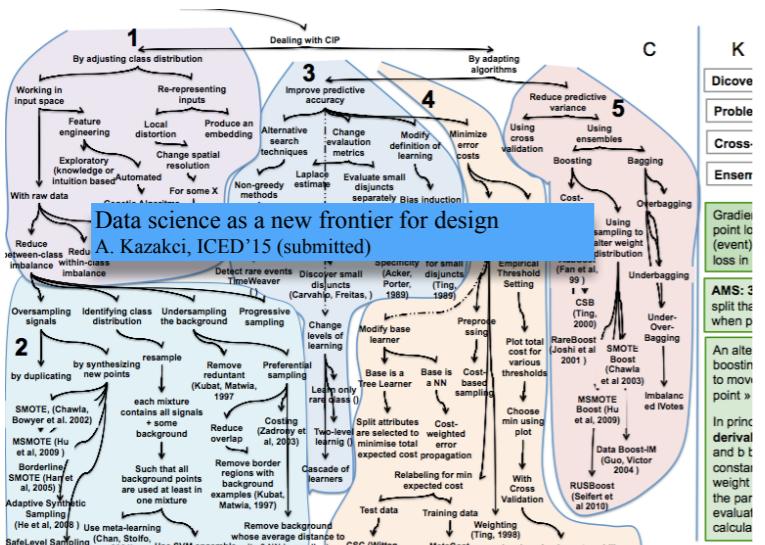
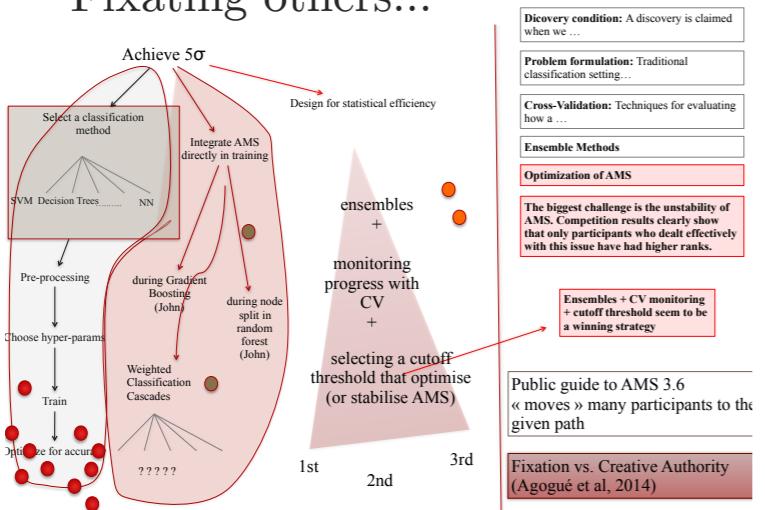


Data challenges are hard to analyze

Analysis of design strategies



Fixating others...



- Available data for HiggsML
 - Forums → 136 topics, 1400+ posts
 - Documentation
 - Participants' blog entries
 - GitHub codes
- Qualitative interpretation combined with C-K modelling of participants' strategies

How do you put a crowd under a microscope?



The research instrument

RAMP - Rapid Analytics and Model Prototyping

A Collaborative Development Platform for Data Science

my submissions
new submission
leaderboard
log out

Combined score: 0.227
Combined test score: 0.225

Number of air passengers prediction

Instant access to all submitted code - for participants & organizers

Private leaderboard

team	submission	public bag
mouhcine_taoufig	First submition	0.244
robin_vogel	sparse and spurious	0.237
baptiste_roziere	INOPA	0.263
nicolas_legroux	second_try	0.260
yousef_sebiat	last_sub	0.242
yousef_sebiat	HMM_2	0.241
pierre_fournier	firstTryOKLM	0.257
yousef_sebiat	HMM_3	0.241

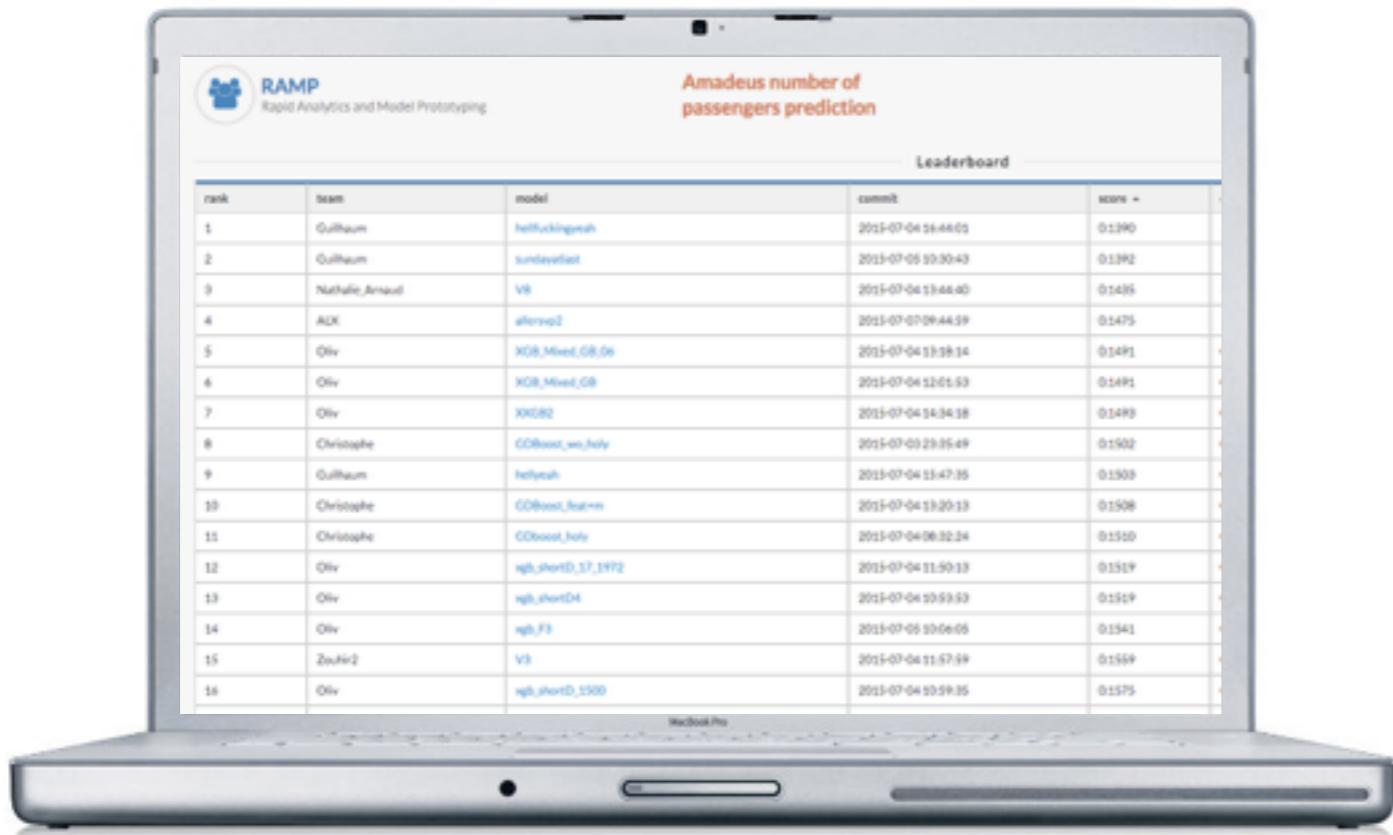
Databoard

Leaderboard > kegl > MF.AB(20;RF(100;5))_d1 > model.py

```
1. from sklearn.ensemble import RandomForestClassifier, AdaBoostClassifier
2. from sklearn.preprocessing import Imputer
3. from sklearn.pipeline import Pipeline
4. from sklearn.base import BaseEstimator
5.
6. class Classifier(BaseEstimator):
7.     def __init__(self):
8.         self.clf = Pipeline([('imputer', Imputer(strategy='most_frequent')),
9.                             ('rf', AdaBoostClassifier(base_estimator=RandomForestClassifier(max_depth=5,
n_estimators=100),
n_estimators=20))])
10.
11.
12.     def fit(self, X, y):
13.         self.clf.fit(X, y)
14.
15.     def predict(self, X):
16.         return self.clf.predict(X)
17.
18.     def predict_proba(self, X):
19.         return self.clf.predict_proba(X)
20.
```

RAMP - Rapid Analytics and Model Prototyping

A Collaborative Development Platform for Data Science



RAMP allows us to collect data on the data science model development process

- 1** We prepare a 'starting kit'
- 2** Continuous access to code:
 - Participants can analyse and build on every submission
 - Organizers can follow real-time what's happening - and react
- 3** Submissions are trained and performances are displayed
- 4** Users actions and interactions are recorded
- 5** Main Output: Dozens of predictive models and performance benchmark

Collecting data with RAMP



RAMP
Rapid Analytics and Model Prototyping

Number of air passengers prediction

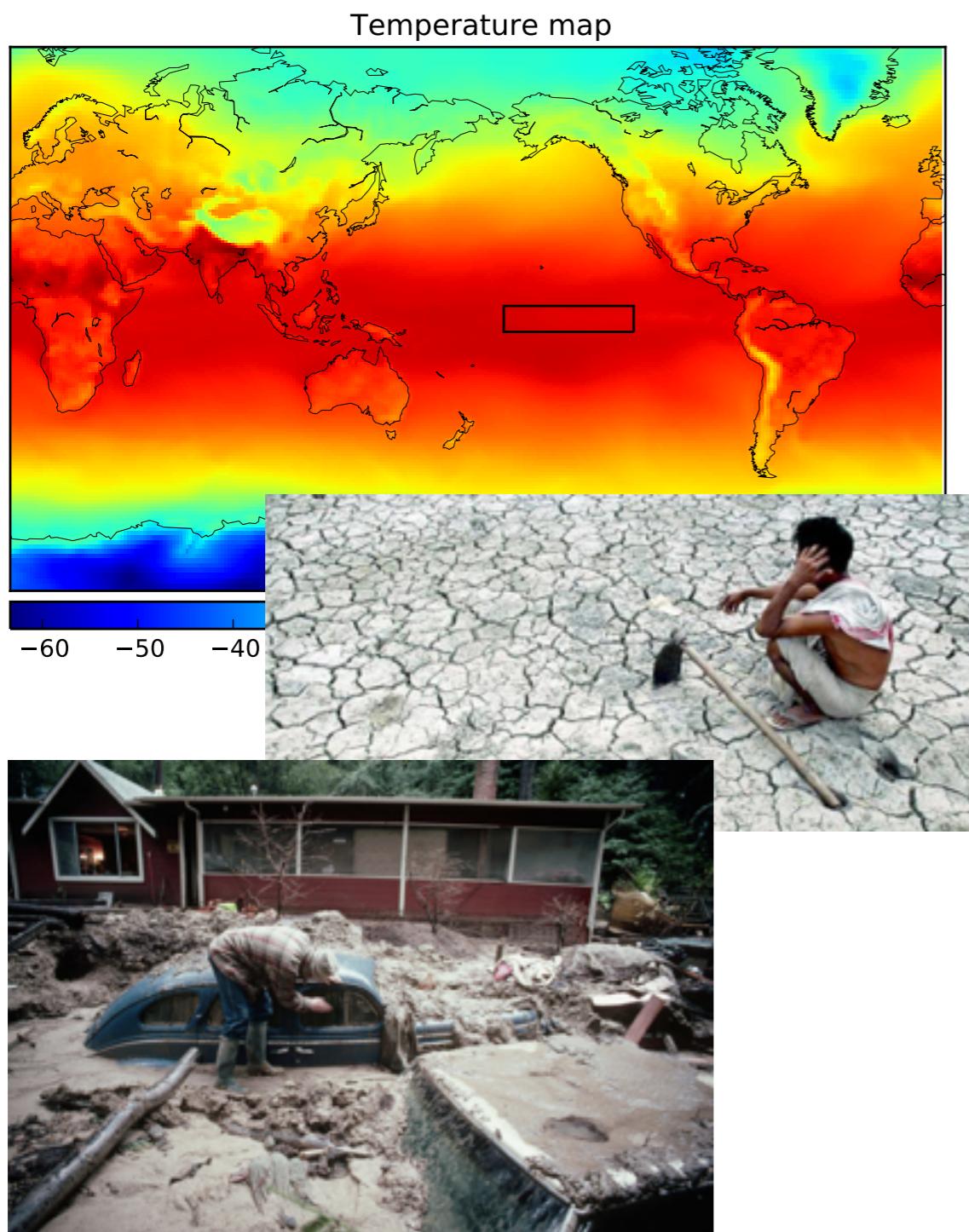
my submissions
new submission
leaderboard
log out

- Number of submissions
 - Frequency of submissions
 - Timing of submissions
 - User interactions
 - Performance of submissions
 - Submitted code



Some applications, preliminary
observations & findings

Climatology



Time Series Based Event Prediction on Geo-tagged data

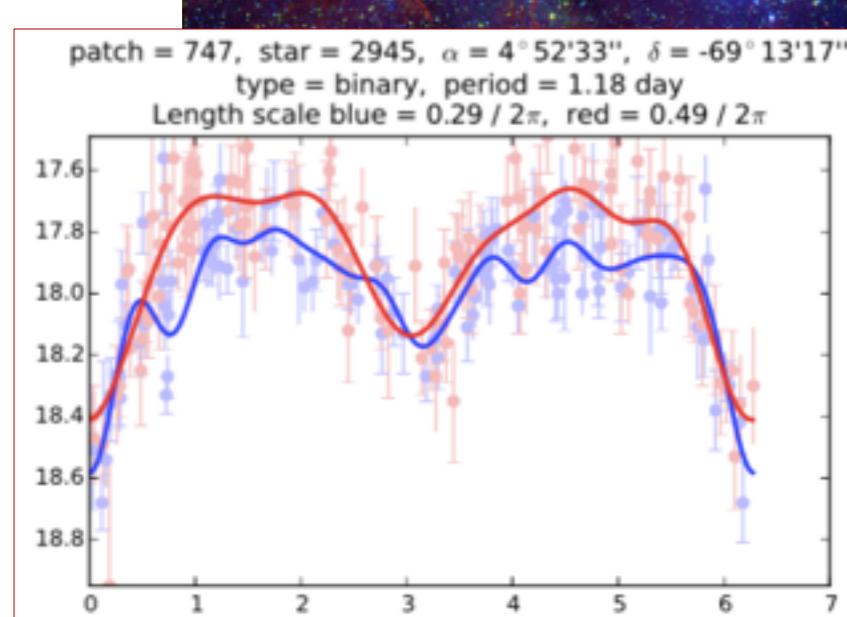
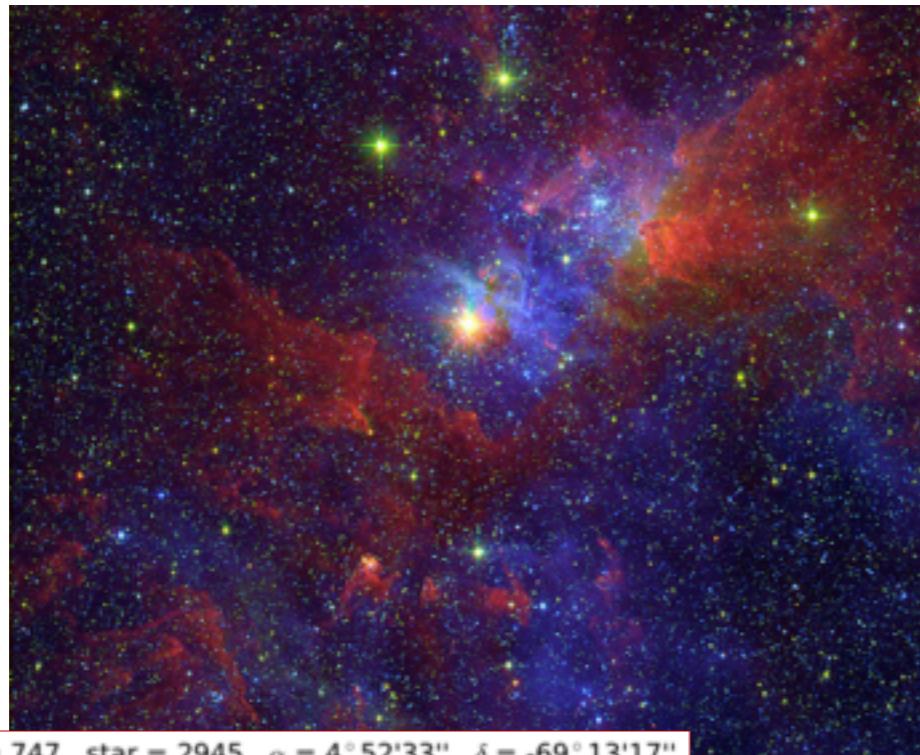
George Washington University
George Mason University

El Niño Prediction

- Temperature data
- Geo-tagged time series
- Prediction: 6 months ahead

Two workshops: Improvement in RMSE score: from 0.90 to 0.43

Astrophysics



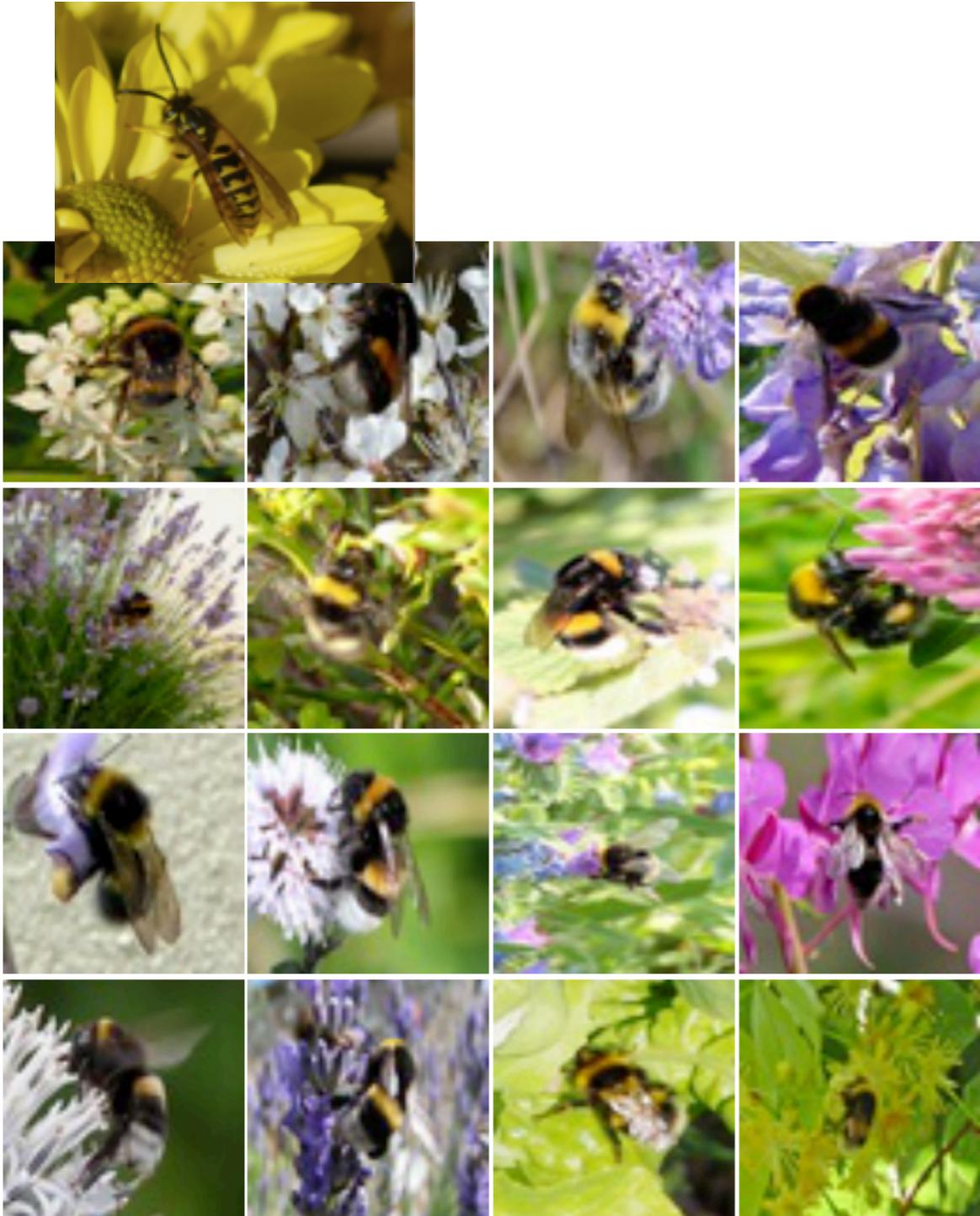
Classification of variable stars

Marc Monier (LAL),
Gilles Faÿ (Centrale-Supelec)

Light curves (luminosity vs time profiles)
- Static features
- Functional data

One day workshop: Accuracy improvement: %89 to %96

Ecology



Finding & Classifying Insects
from Image Data

Paris Museum of Natural History,
SPIPOLL.org, NVDIA,
Université de Champagne-Ardenne ROMEO HPC Center

Pollinating Insects
- Image data (20K images)
- 18 types of insects
- Deep neural net models

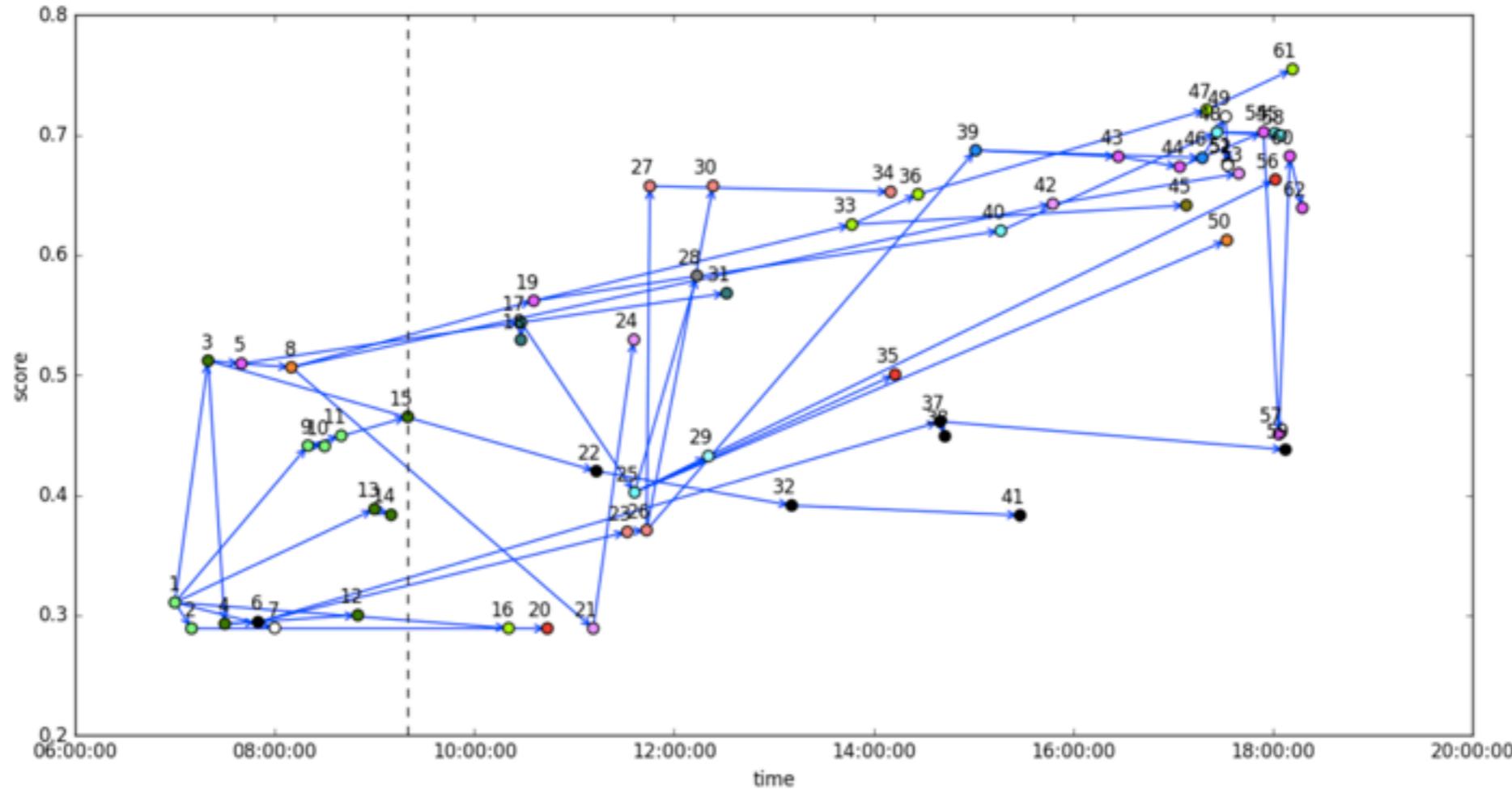
One day event:
Improvement in prediction
accuracy: from 0.31 to 0.70

Cancer Research



Epidemium est un programme de recherche scientifique participatif et ouvert dédié à la compréhension du cancer grâce aux Big Data, qui se concrétisera sous la forme d'un *data challenge*.

A graph of model similarities



Some observations

- **Steady progression.** They have built systematically on a submission they previously created, without being influenced by the others. Their performance may go either up (constantly) or down (constantly).
- **Breakdowns or jumps.** There are other groups, where the performance increased or decreased strongly from one modification submission to the next). There may be some robustness/vulnerability issue with their approach - to be further investigated.
- **Successful expansions.** An important “break” has happened at 12:00. This corresponds to “cropping” idea. Strangely, two very similar submissions (small distance) have been submitted at the same time - one of them did not improve the score at all (around 0.35, while the leader was around 0.55), whereas the other improved considerably (0.65).
- Currently, **we see no dependency between this break and the winning solution.** This might be related to the way we have measured the code similarity.

“Note that, following the RAMP approach, this model is the result of **a succession of small improvements made on top of other participants’ contributions**. We did not reach a prediction score of 0.71 in one shot, but after applying several tricks and manually tuning some parameters.”

Heuritech,
Winner of Insect Challenge
Blog entry

How to compare design concepts - represented as code?

RAMP
Rapid Analytics and Model Prototyping

Number of air passengers prediction

my submissions new submission leaderboard log out

Leaderboard > hien-long-bruno_ly > GBYTES > feature_extractor.py

```
1. import pandas as pd
2. import os
3.
4. class FeatureExtractor(object):
5.     def __init__(self):
6.         pass
7.
8.     def fit(self, X_df, y_array):
9.         pass
10.
11.    def transform(self, X_df):
12.        X_encoded = X_df
13.        path = os.path.dirname(__file__)
14.        data_weather = pd.read_csv(os.path.join(path,
15. "external_data.csv"))
16.        X_weather = data_weather[['Date', 'Airport', 'Max
TemperatureC']]
17.        X_weather = X_weather.rename(
18.            columns={'Date': 'DateOfDeparture', 'Airport': 'Arrival'})
19.        X_encoded = X_encoded.set_index(['DateOfDeparture', 'Arrival'])
20.        X_weather = X_weather.set_index(['DateOfDeparture', 'Arrival'])
21.        X_encoded = X_encoded.join(X_weather).reset_index()
22.
23.        X_encoded = X_encoded.join(pd.get_dummies(
24.            X_encoded['Departure'], prefix='d'))
25.        X_encoded = X_encoded.join(
26.            pd.get_dummies(X_encoded['Arrival'], prefix='a'))
27.        X_encoded = X_encoded.drop('Departure', axis=1)
28.        X_encoded = X_encoded.drop('Arrival', axis=1)
29.
30.        X_encoded = X_encoded.drop('DateOfDeparture', axis=1)
31.        X_array = X_encoded.values
32.        return X_array
```

~ ?

RAMP
Rapid Analytics and Model Prototyping

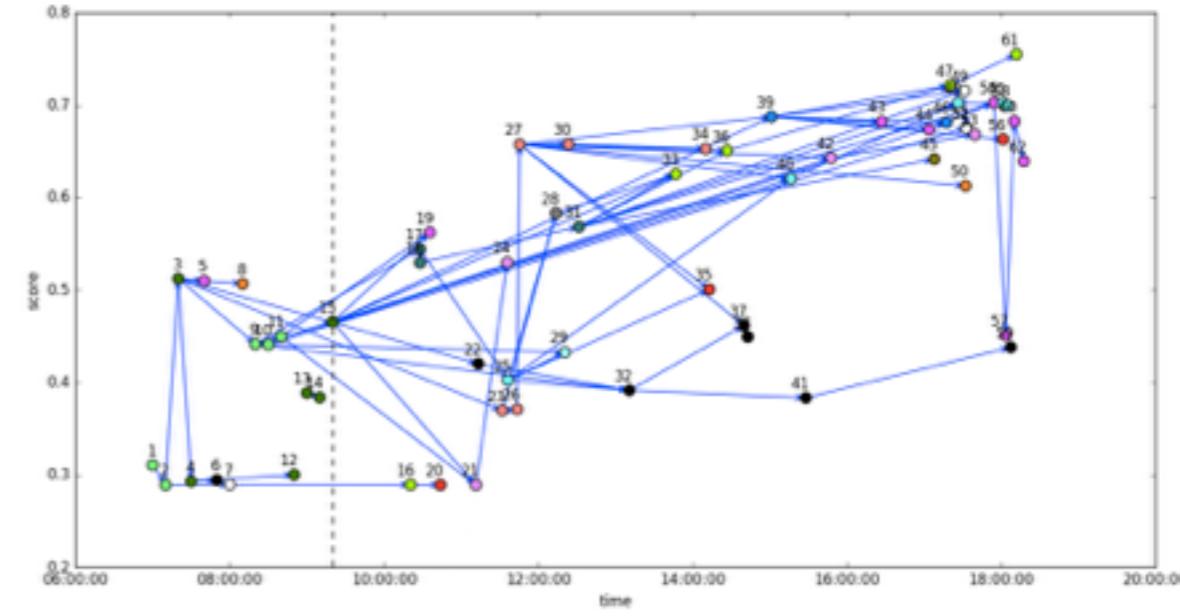
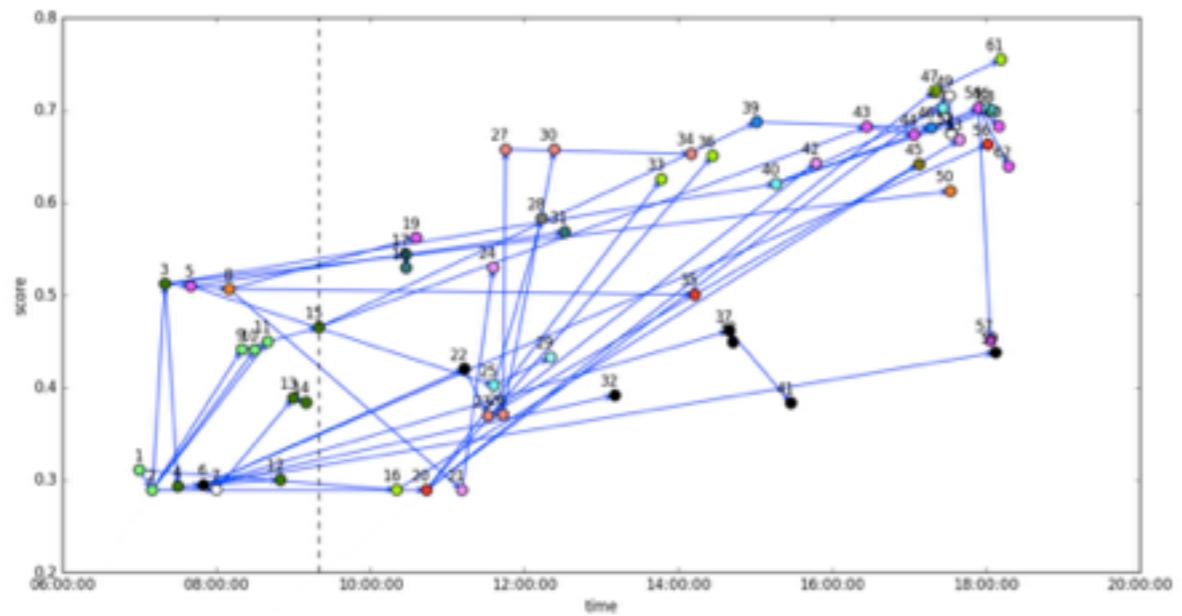
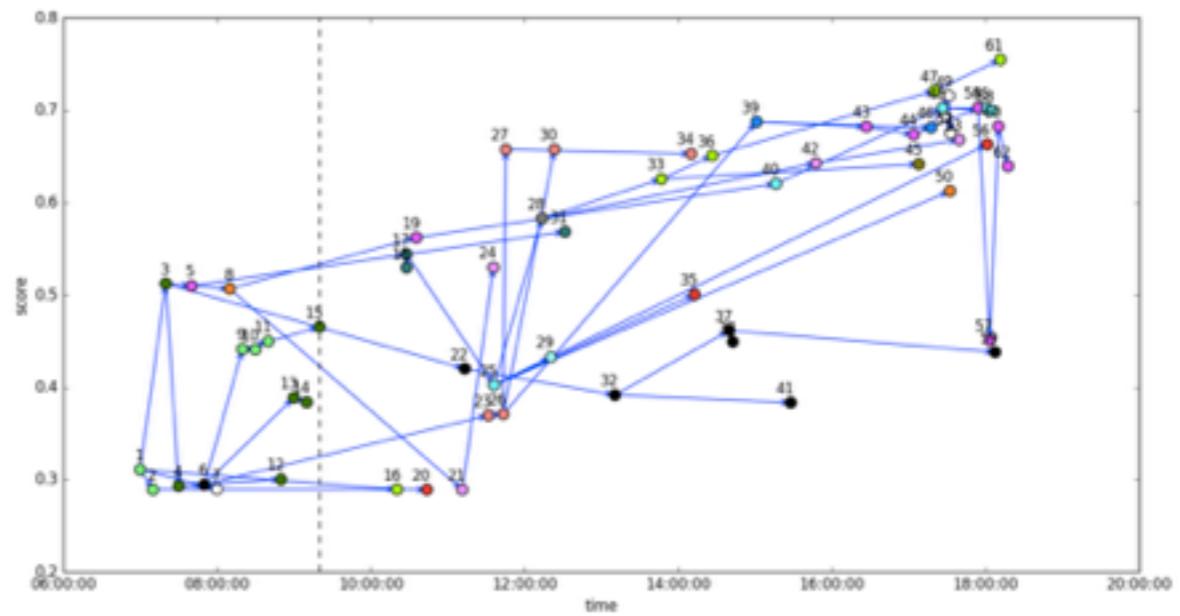
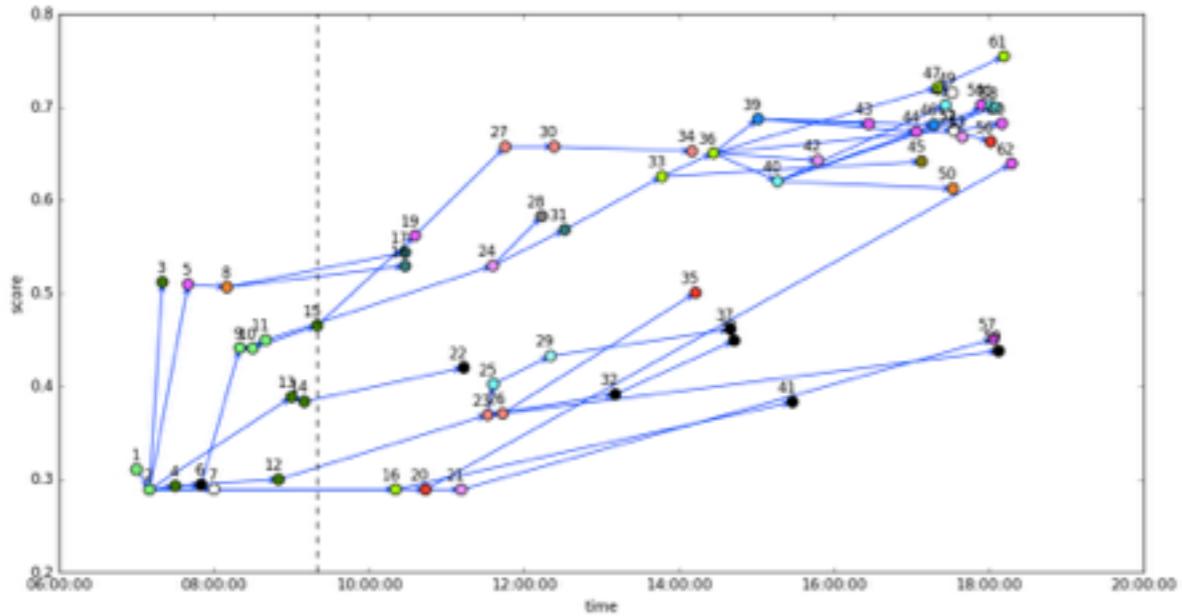
Number of air passengers prediction

my submissions new submission leaderboard log out

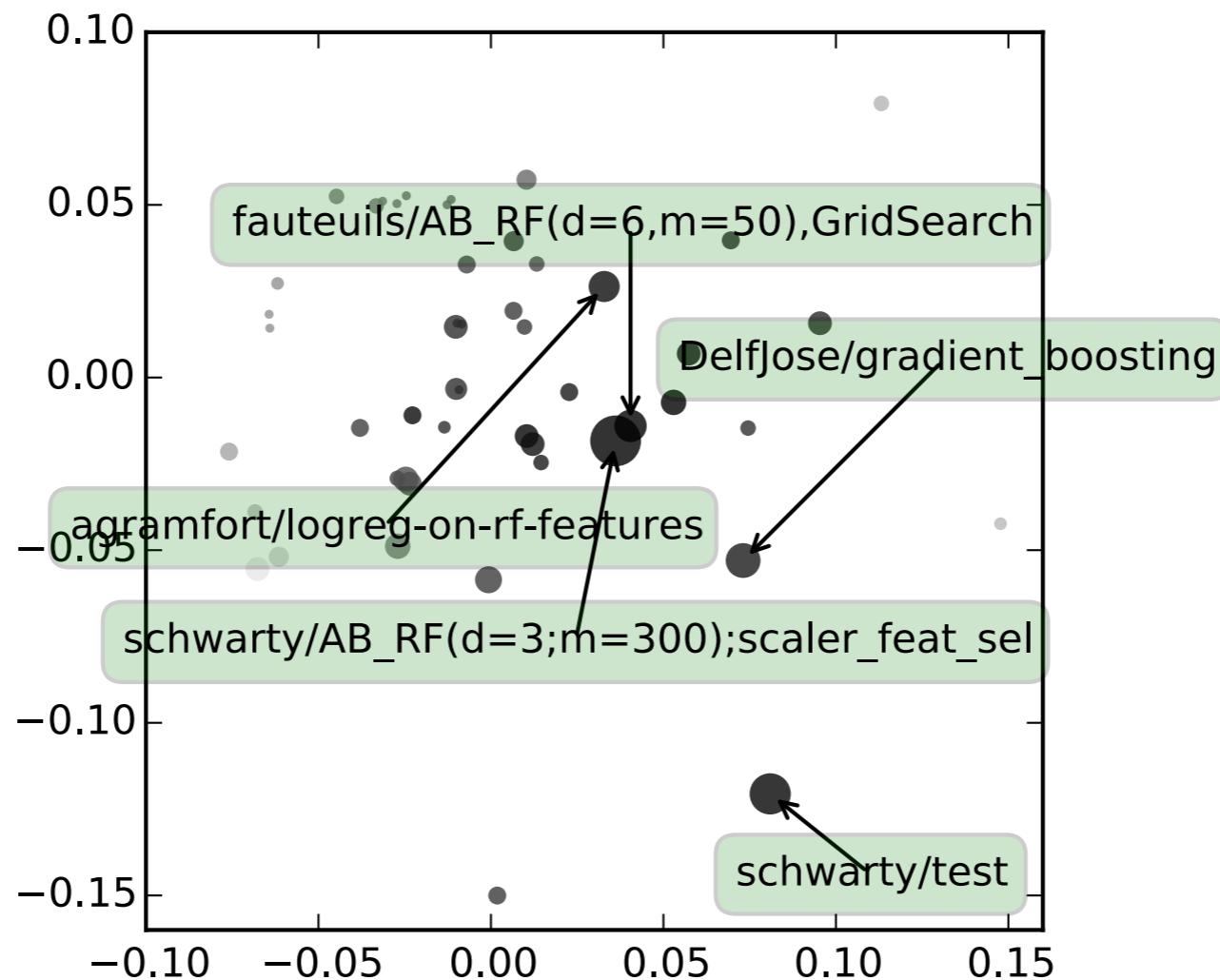
Leaderboard > djabbz > NdioumJob > feature_extractor.py

```
1. import pandas as pd
2. import os
3.
4. def get_binary_var(data, prefix=None, prefix_sep='_'):
5.     from pandas.tools.merge import concat
6.     import numpy as np
7.
8.     data = data.astype('category')
9.     length = data.cat.categories.size
10.    length = int(np.ceil(np.log2(length)))
11.    code_string = data.cat.codes.map(lambda x:
12.        '{0:b}'.format(x).rjust(length, '0'))
13.
14.    codes = []
15.    for i in range(length):
16.        codes.append(code_string.str[i])
17.
18.    if levels = data.cat.categories
19.        if prefix is None:
20.            prefix = data.name
21.            col_name = ('{0}{1}'*(prefix, prefix_sep, 1) for l in range(0,
length + 1))
22.            df = concat(codes, axis=1)
23.            df.columns = col_name
24.            return df.astype(np.int8)
25.    class FeatureExtractor(object):
26.        def __init__(self):
27.            pass
28.
29.        def fit(self, X_df, y_array):
30.            pass
31.
32.        def transform(self, X_df):
33.
34.            data_encoded = X_df
35.
36.            state_pop = pd.DataFrame({('ATL', 9.9),
37.                ('BOS', 6.6),
38.                ('CLT', 9.7),
39.                ('DEN', 5.2),
40.                ('DFW', 26.6),
41.                ('EWR', 9.9),
42.                ('IAD', 8.8),
43.                ('JFK', 26.6),
44.                ('LAX', 19.5),
45.                ('ORD', 2.7),
46.                ('PDX', 38.1),
47.                ('SFO', 19.5),
48.                ('SEA', 19.3),
49.                ('SLC', 19.3),
50.                ('SFO', 5.4),
51.                ('TIA', 12.9),
52.                ('PHL', 12.8),
53.                ('FLL', 6.5),
54.                ('SEA', 6.9),
55.                ('IPO', 38.1), columns=['airports'],
56.                'pop'))
57.
58.            data_encoded = data_encoded.merge(state_pop, how='outer',
59.                left_on='Departure', right_on='airports', axis=1)
59.            data_encoded = data_encoded.merge(state_pop, how='outer',
60.                left_on='Arrival', right_on='airports').drop('airports', axis=1)
61.
62.            data_encoded =
63.            data_encoded.join(get_binary_var(data_encoded['Departure'],
64.                prefix='d'))
64.            data_encoded =
65.            data_encoded.join(get_binary_var(data_encoded['Arrival'],
66.                prefix='a'))
66.            data_encoded = data_encoded.drop('Departure', axis=1)
67.            data_encoded = data_encoded.drop('Arrival', axis=1)
68.
69.            data_encoded['DateOfDeparture'] =
70.                pd.to_datetime(data_encoded['DateOfDeparture'])
71.            data_encoded['year'] = data_encoded['DateOfDeparture'].dt.year
72.            data_encoded['month'] =
73.                data_encoded['DateOfDeparture'].dt.month
74.            data_encoded['day'] = data_encoded['DateOfDeparture'].dt.day
75.            data_encoded['weekday'] =
76.                data_encoded['DateOfDeparture'].dt.weekday
77.            data_encoded['n_days'] =
78.                data_encoded['DateOfDeparture'].apply(lambda date: (date -
79.                    pd.to_datetime('1970-01-01')).days)
80.
81.            data_encoded = data_encoded.drop('DateOfDeparture', axis=1)
82.
83.            data_encoded =
84.            data_encoded.join(get_binary_var(data_encoded['year'],
85.                prefix='y'))
86.            data_encoded =
87.            data_encoded.join(get_binary_var(data_encoded['month'],
88.                prefix='m'))
```

Various distance measures



Comparing performance profiles: Promoting novelty search



2D projection (MDS) of model's prediction profiles

- Greyness: Model's raw score
- Size: model's contribution
- Position: similarity/dissimilarity in predictions

We Found (to be validated by further studies):

- Gravitation: following a given submission, others are hovering around the same coordinates, by incremental adjustments
- Repulsion: new submission using out-of-the-box code to explore the white space (no previous close-by submissions exist)
- Hybridation: opportunistic integration of previous submissions, involving/inspired by at least two different source of code.

In progress:

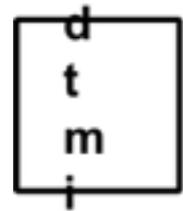
- Monitoring & Modelling “contribution”
(Pierre Fleckinger, Economic Agents & Incentive Theory)
- Pushing towards “Novelty Search”
(Jean-Bastiste Mouret, Novelty Search)
- Controlled experiments

RAMP platform

- RAMP platform is meant to be a free tool for researchers and students; this opens up new perspectives (pedagogy & research) and hopefully brings closer different communities

Thank you

Akin Kazakci, Mines ParisTech
akin.kazakci@mines-paristech.fr



design theory
and methods
for innovation

