

université
PARIS-SACLAY

 Paris-Saclay
Center for Data Science

Open Software Initiative (OSI)

Alexandre Gramfort

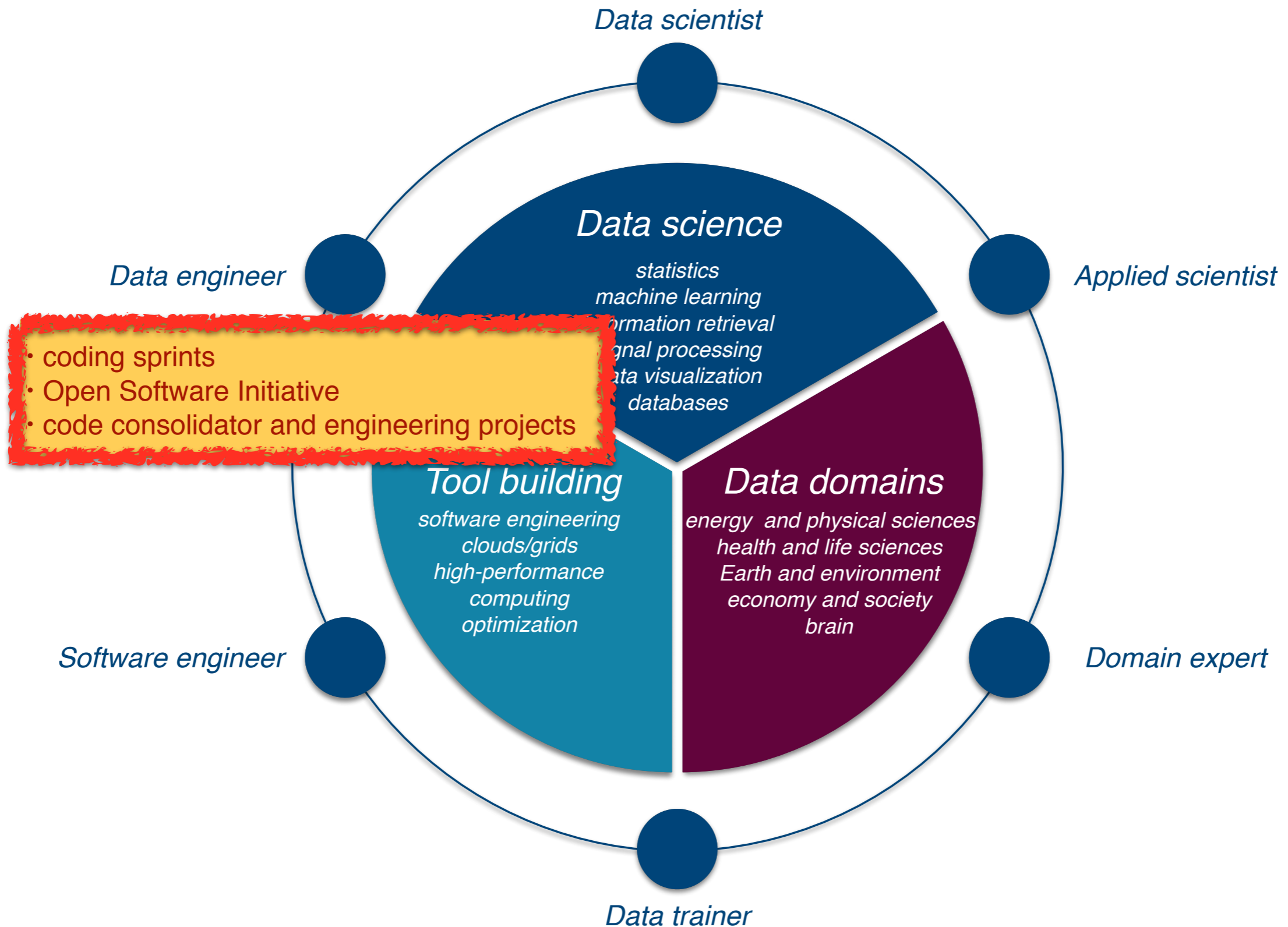
alexandre.gramfort@telecom-paristech.fr

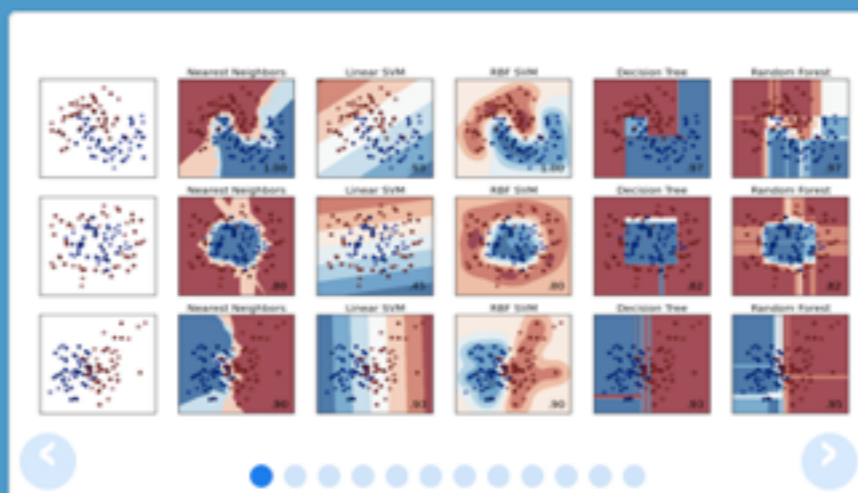
Assistant Professor

LTCI, CNRS, Télécom ParisTech, Université Paris-Saclay



TOOLS IN THE DATA SCIENCE LANDSCAPE





scikit-learn

Machine Learning in Python

- Simple and efficient tools for data mining and data analysis
- Accessible to everybody, and reusable in various contexts
- Built on NumPy, SciPy, and matplotlib
- Open source, commercially usable - BSD license

Classification

Identifying to which set of categories a new observation belong to.

Applications: Spam detection, Image recognition.

Algorithms: *SVM, nearest neighbors, random forest, ...* — Examples

Regression

Predicting a continuous value for a new example.

Applications: Drug response, Stock prices.

Algorithms: *SVR, ridge regression, Lasso, ...* — Examples

Clustering

Automatic grouping of similar objects into sets.

Applications: Customer segmentation, Grouping experiment outcomes

Algorithms: *k-Means, spectral clustering, mean-shift, ...* — Examples

Dimensionality reduction

Reducing the number of random variables to consider.

Applications: Visualization, Increased efficiency

Algorithms: *PCA, Isomap, non-negative matrix factorization.* — Examples

Model selection

Comparing, validating and choosing parameters and models.

Goal: Improved accuracy via parameter tuning

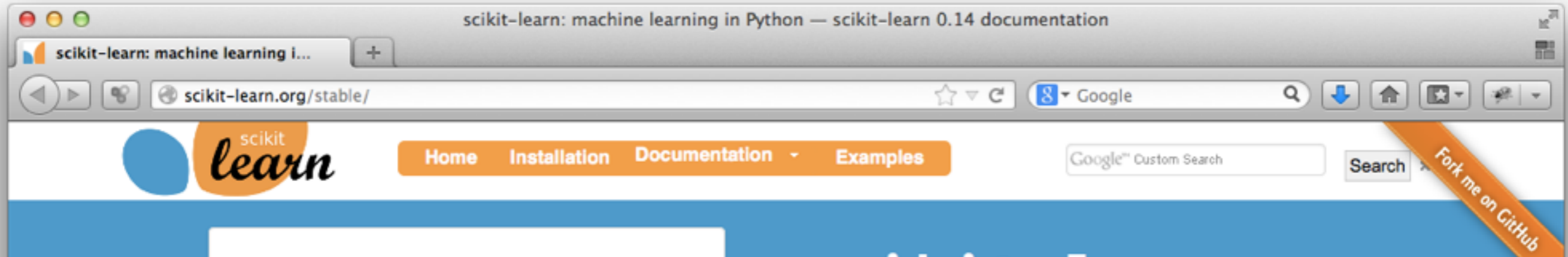
Modules: *grid search, cross validation, metrics.* — Examples

Preprocessing

Feature extraction and normalization.

Application: Transforming input data such as text for use with machine learning algorithms.

Modules: *preprocessing, feature extraction.* — Examples



In a Nutshell, scikit learn...

Started in 2010 at:

université
PARIS-SACLAY

... has had **20,181 commits** made by **650 contributors** representing **179,355 lines of code**

... is **mostly written in Python** with a **well-commented source code**

... has a **well established, mature codebase** maintained by a **very large development team** with **stable Y-O-Y commits**

... took an estimated **47 years of effort** (COCOMO model) starting with its **first commit in January, 2010** ending with its **most recent commit 3 days ago**

source: <https://www.openhub.net/p/scikit-learn>

Classification

Identifying to which class an observation belongs to

Applications: Spam recognition.

Algorithms: SVM, random forest, ...

Dimensionality

Reducing the number of features to consider.

Applications: Visual efficiency

Algorithms: PCA, Isomap, matrix factorization.



In a Nutshell, scikit learn...

Started in 2010 at:



... has had 20,181 commits representing 175k lines of code

... is mostly written in Python with a well-maintained code base

... has a well-maintained code base with stable releases

... took an estimated 1.5 years starting with its first release ending with its most recent

Funding:



Paris-Saclay Center for Data Science



Classification

Identifying to which class an observation belongs to. Applications: Spam recognition. Algorithms: SVM, random forest, ...

Dimensionality

Reducing the number of features to consider. Applications: Visual efficiency. Algorithms: PCA, Isomap, matrix factorization.



DATA PUBLICA



EVERNOTE

<http://scikit-learn.org/stable/testimonials/testimonials.html>

PeerIndex

fnrs
LA LIBERTÉ DE CHERCHER

Inria
INVENTORS FOR THE DIGITAL WORLD

TELECOM
ParisTech

DataRobot

lovely

ŷhat

how
about
we...

okc



INFONEA
Comma Soft AG



DATA PUBLICA

- Installed on 1% of debian systems
- 1200 job offers on Stack-Overflow
- Users: 60% academics & 40% industry

d
Writ



BESTOT
MEDIA
GROUP

RANGE SPAN

EVERNOTE

BEYOND SCIKIT-LEARN @ UPSA

(... IN MY WORLD)

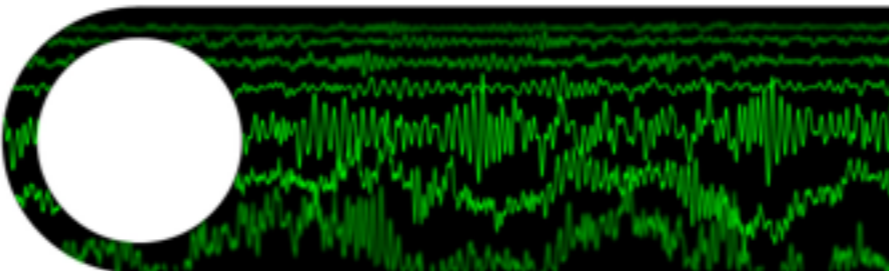
MNE

MEG + EEG ANALYSIS & VISUALIZATION



ni
learn

neo



NiBabel

Access a cacophony of neuro-imaging file formats

and certainly many more...

Objective is to teach students to replicate the scikit-learn model or at least contribute to it



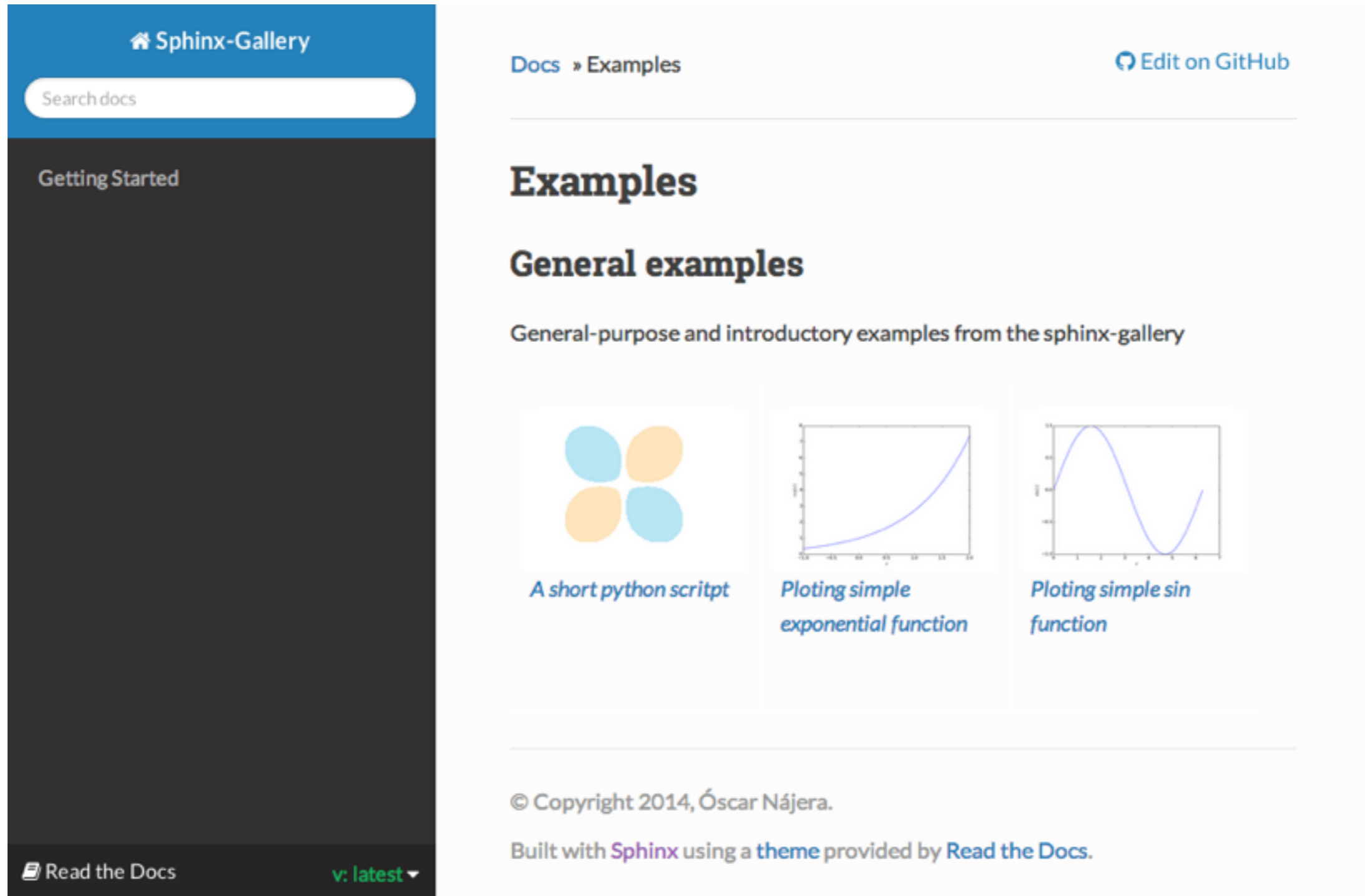
Open Software Initiative @ CDS I

- 6 projects funded (5 missions doctorales + 1 code consolidator 3 months)
 - N. Goix (LTCI) : Scikit-Learn “Outlier/Novelty detection algorithm (IForest, etc.)”
 - M. Cherti (UpSud) : Scikit-Learn “Multivariate Additive Regression Splines (MARS)”
 - R. Brault (IBISC) : OperaLib “Large scale learning with operator valued kernels”
 - O. Najera (UpSud) : NiLearn “Sphinx-Gallery: Facilitate code documentation with generated example gallery”
 - S. Mishra (UpSud) : OpenMEEG/MNE “Integration of OpenMEEG into the MNE software via Python bindings”
 - Lorenzo De Santis (UpSud) : MNE “Port from C to Python dipole fitting solver for source localization”
 - Xin Su (LTCI) : START² “Synthetic Aperture Radar Time series Toolbox”

Open Software Initiative @ CDS I

- 6 projects funded (5 missions doctorales + 1 code consolidator 3 months)
 - N. Goix (LTCI) : Scikit-Learn “Outlier/Novelty detection algorithm (IForest, etc.)”
 - M. Cherti (UpSud) : Scikit-Learn “Multivariate Additive Regression Splines (MARS)”
 - R. Brault (IBISC) : OperaLib “Large scale learning with operator valued kernels”
 - O. Najera (UpSud) : NiLearn “Sphinx-Gallery: Facilitate code documentation with generated example gallery”
 - S. Mishra (UpSud) : OpenMEEG/MNE “Integration of OpenMEEG into the MNE software via Python bindings”
 - Lorenzo De Santis (UpSud) : MNE “Port from C to Python dipole fitting solver for source localization”
 - Xin Su (LTCI) : START² “Synthetic Aperture Radar Time series Toolbox”

OSI : O. Najera “Sphinx-Gallery”



The screenshot displays the Sphinx-Gallery website interface. On the left is a dark sidebar with a blue header containing the site name and a search bar. The main content area is white and features a breadcrumb trail, a GitHub edit link, and a list of example categories. Three example cards are visible, each with a thumbnail image and a title. The footer contains copyright information and build details.

Sphinx-Gallery

Search docs


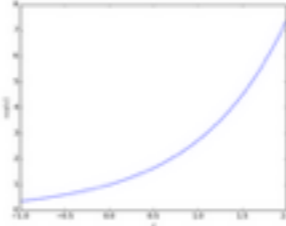
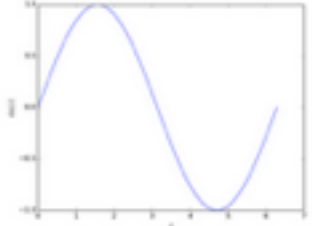
Getting Started

Docs » Examples [Edit on GitHub](#)

Examples

General examples

General-purpose and introductory examples from the sphinx-gallery

- 
A short python script
- 
Ploting simple exponential function
- 
Ploting simple sin function

© Copyright 2014, Óscar Najera.
Built with [Sphinx](#) using a [theme](#) provided by [Read the Docs](#).

[Read the Docs](#) v: latest

<http://sphinx-gallery.readthedocs.io/en/latest/>

<https://github.com/sphinx-gallery/sphinx-gallery>

[Getting Started to Sphinx-Gallery](#)[Advanced Configuration](#)[Sphinx-Gallery syntax](#)[Sphinx-Gallery API Reference](#)[Gallery of Examples](#)[Secondary gallery](#)[Change Log](#)

Your vote matters! Find your local polling place and get out to vote on November 8th

Welcome to Sphinx-Gallery's documentation!

A [Sphinx](#) extension that builds an HTML gallery of examples from any set of Python scripts.

It is extracted from the scikit-learn project and aims to be an independent general purpose extension.

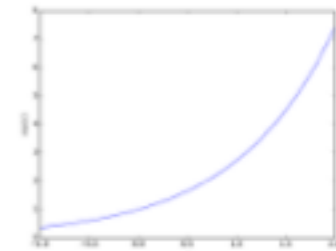
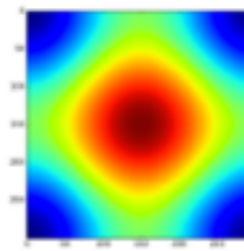
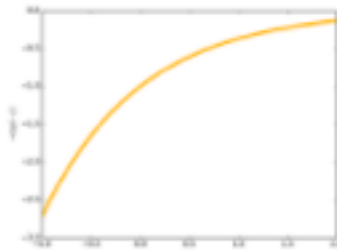
Why Sphinx-Gallery?

- Simple examples that run out of the box are the best way to learn a library
- Pleasing, organized, visual layouts
- Links, searching, backlinks throughout examples and documentation

What does it look like?

Here is an example gallery generated from a few Python scripts.

Examples using `numpy.linspace`

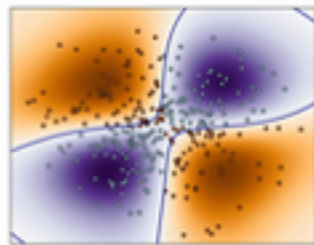


OSI : O. Najera “Sphinx-Gallery”

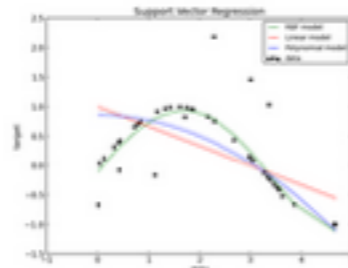
http://scikit-learn.org/stable/auto_examples/index.html

Support Vector Machines

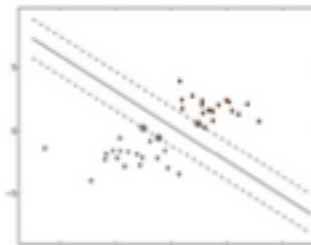
Examples concerning the `sklearn.svm` package.



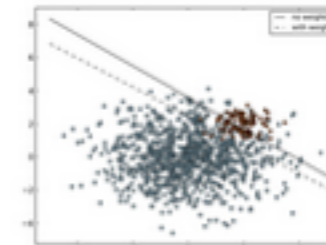
Non-linear SVM



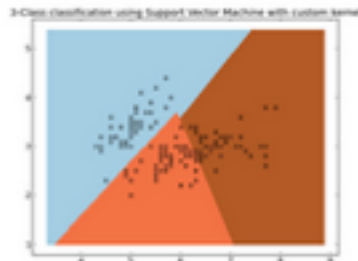
Support Vector Regression (SVR) using linear and non-linear kernels



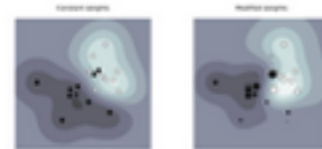
SVM: Maximum margin separating hyperplane



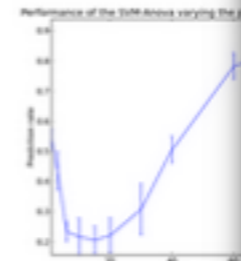
SVM: Separating hyperplane for



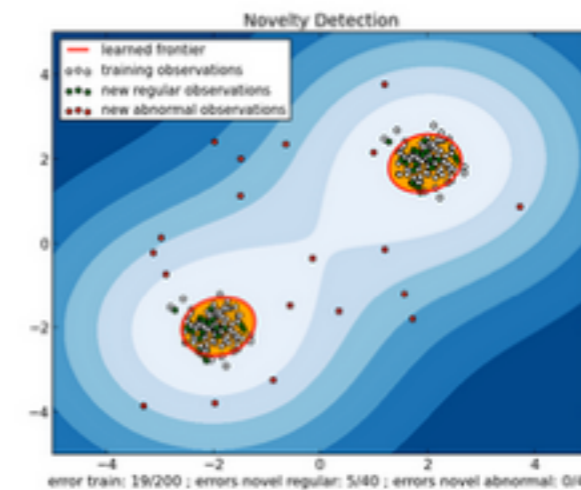
SVM with custom kernel



SVM: Weighted samples



SVM-Anova: SVM univariate feature selection



One-class SVM with non-linear kernel (RBF)

An example using a one-class SVM for novelty detection.

OSI : O. Najera “Sphinx-Gallery”

http://scikit-image.org/docs/dev/auto_examples/



scikit-image
image processing in python

[Home](#)

[Download](#)

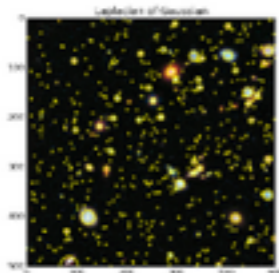
[Gallery](#)

[Documentation](#)

[Source](#)

General examples

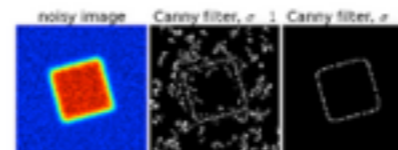
General-purpose and introductory examples for the scikit.



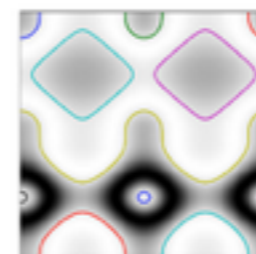
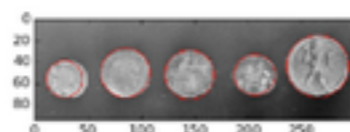
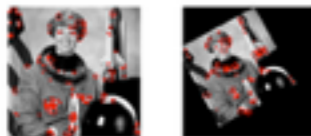
Blob Detection



BRIEF binary descriptor



Canny edge detector



Navigation

[Documentation Home](#)

Previous topic

[License](#)

Next topic

[Blob Detection](#)

Contents

[General examples](#)

[Longer examples and demonstrations](#)

Versions

[skimage dev](#)

[skimage 0.10.x](#)

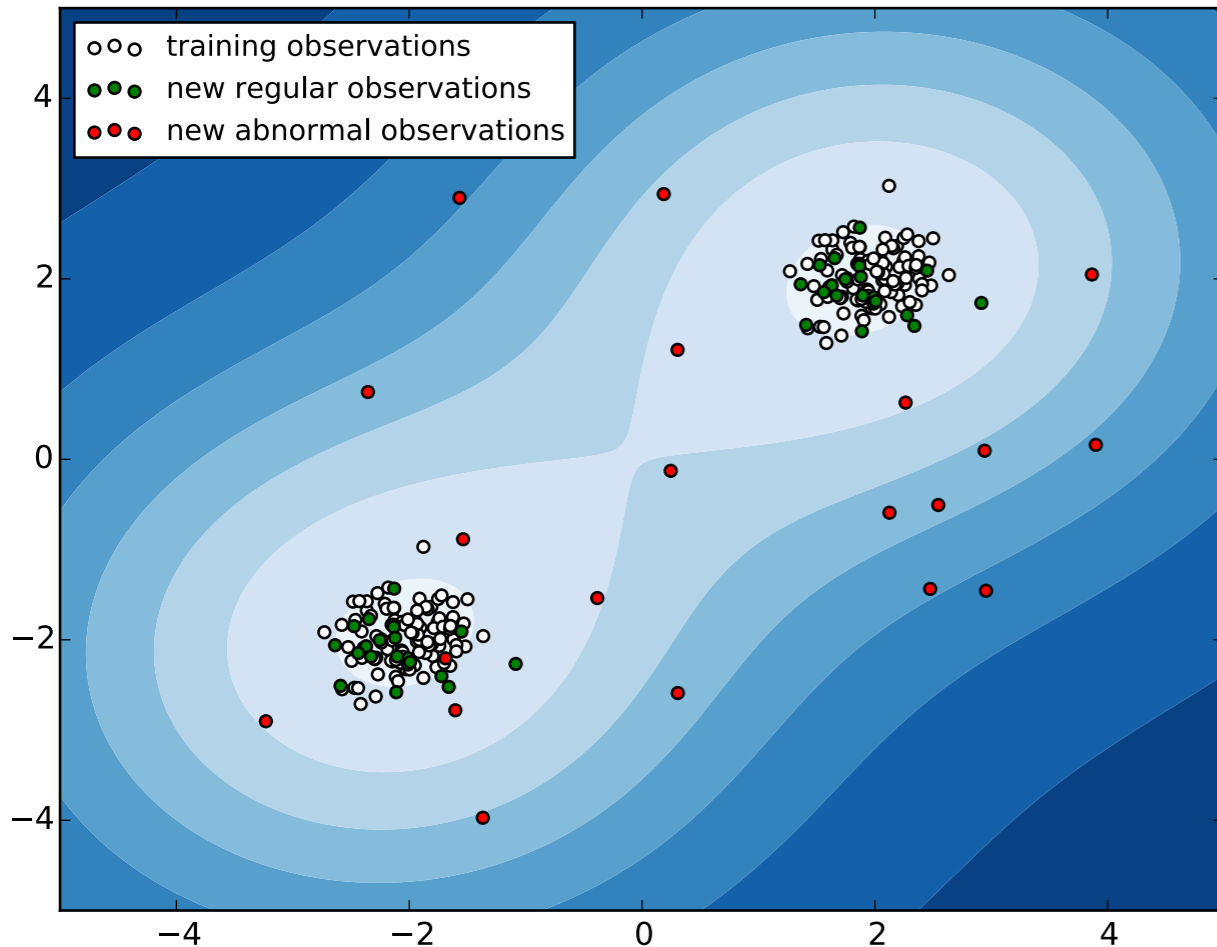
[skimage 0.9.x](#)

Open Software Initiative @ CDS I

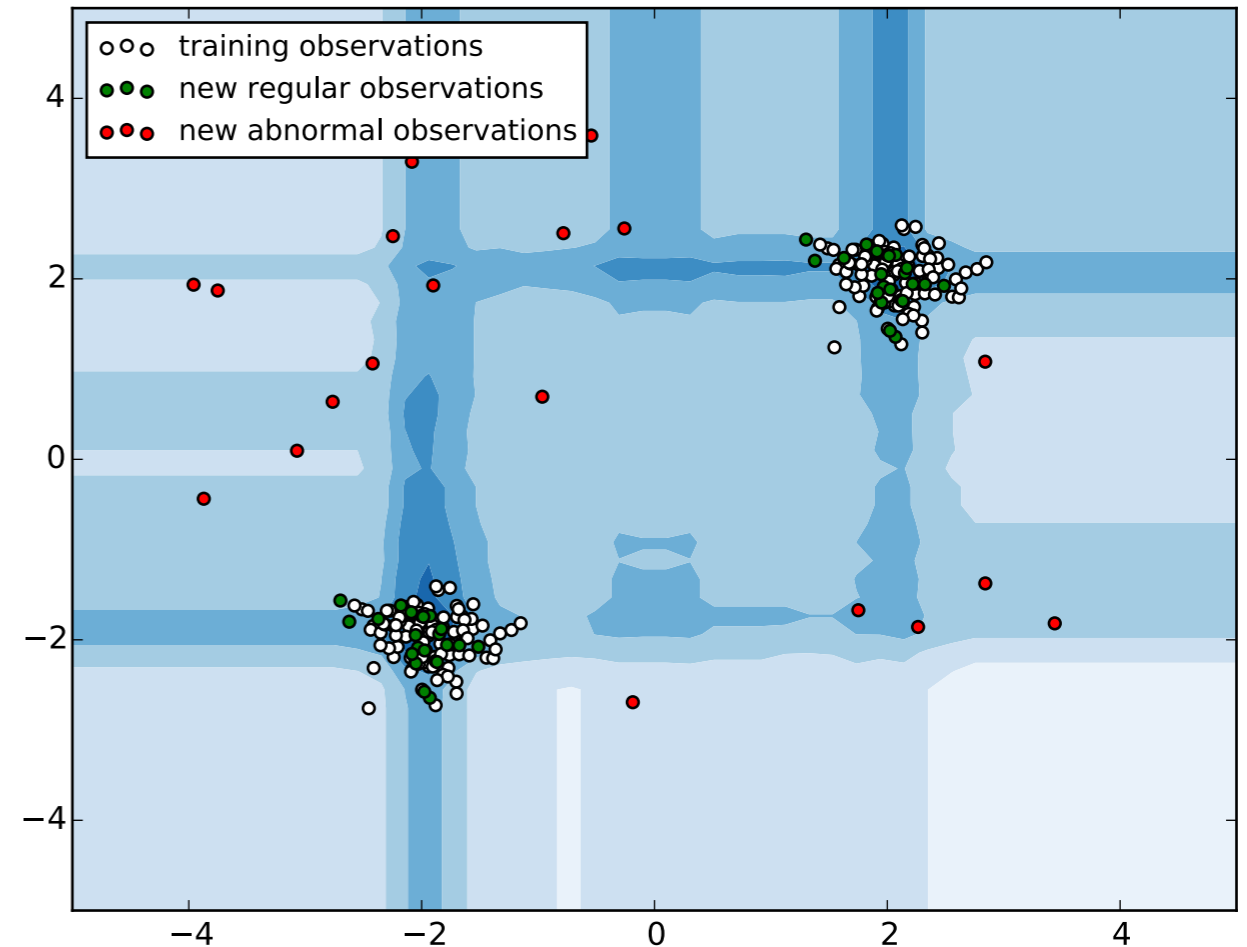
- 6 projects funded (5 missions doctorales + 1 code consolidator 3 months)
 - N. Goix (LTCl) : Scikit-Learn “Outlier/Novelty detection algorithm (IForest, etc.)”
 - M. Cherti (UpSud) : Scikit-Learn “Multivariate Additive Regression Splines (MARS)”
 - R. Brault (IBISC) : OperaLib “Large scale learning with operator valued kernels”
 - O. Najera (UpSud) : NiLearn “Sphinx-Gallery: Facilitate code documentation with generated example gallery”
 - S. Mishra (UpSud) : OpenMEEG/MNE “Integration of OpenMEEG into the MNE software via Python bindings”
 - Lorenzo De Santis (UpSud) : MNE “Port from C to Python dipole fitting solver for source localization”
 - Xin Su (LTCl) : START² “Synthetic Aperture Radar Time series Toolbox”

OSI : N. Goix “Outlier detection in Scikit-Learn”

One class SVM



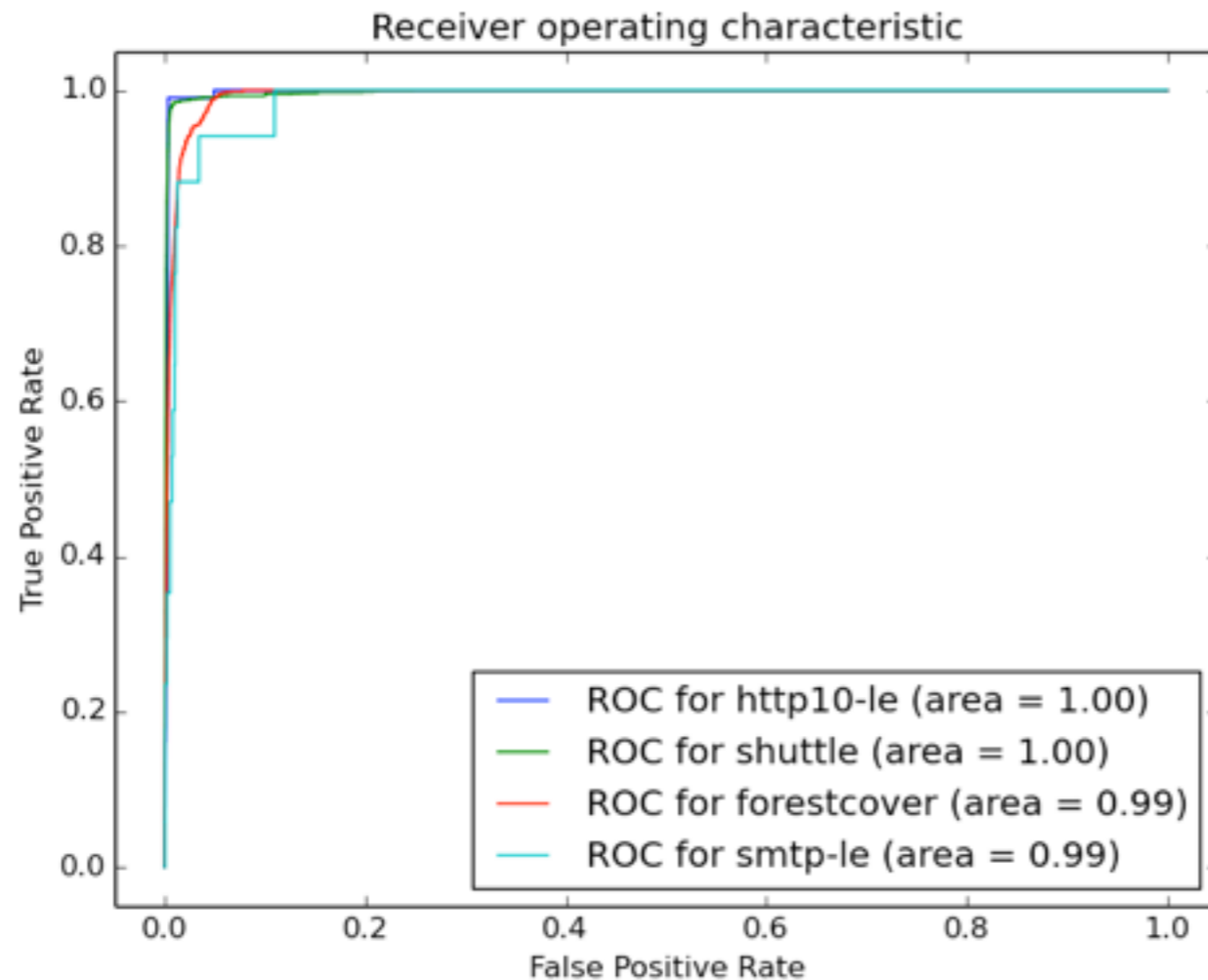
Outlier Detection



<http://scikit-learn.org/stable/modules/generated/sklearn.ensemble.IsolationForest.html>

F.T. Liu, Kai Ming Ting, and Zhi-Hua Zhou. Isolation forest. In Data Mining, 2008. ICDM '08. Eighth IEEE International Conference on, pages 413--422, Dec 2008.

OSI : N. Goix “Outlier detection in Scikit-Learn”



Result on
KDD 99 dataset

<http://scikit-learn.org/stable/modules/generated/sklearn.ensemble.IsolationForest.html>

F.T. Liu, Kai Ming Ting, and Zhi-Hua Zhou. Isolation forest. In Data Mining, 2008. ICDM '08. Eighth IEEE International Conference on, pages 413--422, Dec 2008.

How to apply?

- Process:
 - Software project proposes eventually a topic. List from 2015:
 - <https://github.com/paris-saclay-cds/open-software-initiative/wiki/Projects>
 - Student applies by writing a project proposal of 2 pages with a roadmap
 - Proposal is refined by iterating between mentor and student

Coding sprints

MNE (<http://martinos.org/mne>) and NiLearn (<http://nilearn.github.io/>)

1 week sprint in June 2014 at Telecom ParisTech
15 attendees (from Saclay, UK, Finland, Austria, Germany, USA)



CONCLUSIONS

- OSI:
 - offers great training for students (team development, collaborative work)
 - offers visibility for the students !
 - teaches students how to write code that others can read so good code !
 - makes students discover a new scientific field