



PyEMMA Package Overview and Software Development

Martin K. Scherer
Free University Berlin

February 17, 2019



Software overview and design patterns

- Python

- Anaconda stack

- Package overview

 - Coordinates package

 - MSM package

PyEMMA Development

- Principles

- Processes

 - GitHub

 - Continous Integration Services

- Collaboration



Software overview and design patterns

- Python

- Anaconda stack

- Package overview

 - Coordinates package

 - MSM package

PyEMMA Development

- Principles

- Processes

 - GitHub

 - Continuous Integration Services

- Collaboration

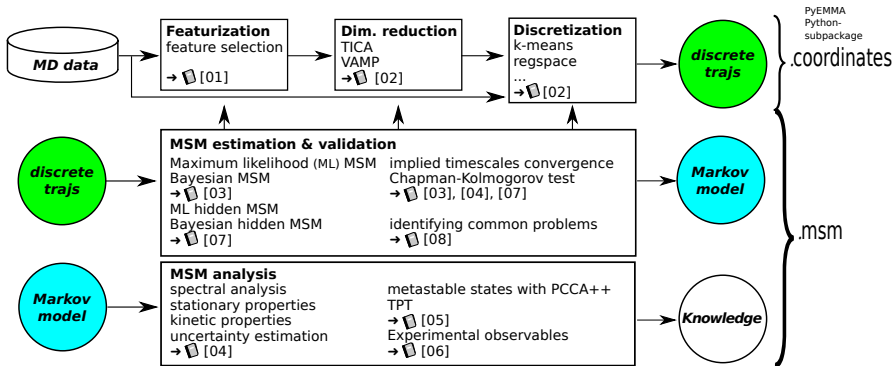


- ▶ Easy to use core libraries (eg. NumPy, SciPy, Pandas, Jupyter, Matplotlib, . . .)
- ▶ Scientific software for MD, data science, biology, chemistry . . .
- ▶ Easy to learn general purpose language
- ▶ Quick prototyping
- ▶ Glue together software written in faster languages (eg. C/C++, Fortran)

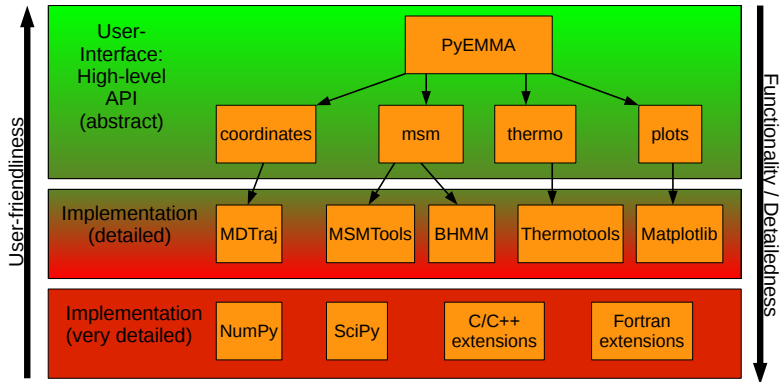


- ▶ Anaconda is a (Python-based) software stack built for all three major platforms (Linux, OSX, Windows)
- ▶ Easy installation and upgrading, no need to compile anything yourself.
- ▶ Different software channels for different purposes (eg. Omnia [MD], BioConda [Bioinformatics], . . .)
- ▶ Automatic handling of dependencies (conflict checking)
- ▶ Possibility to create isolated work environments (separate package versions etc.)

From MD data to Knowledge



Package hierarchy - abstracting detailedness



- ▶ Streaming data pattern
- ▶ Avoid the need of dumping intermediate results to disk
- ▶ Support for multiple data formats
- ▶ Random access possible (either simulated or IO efficient)

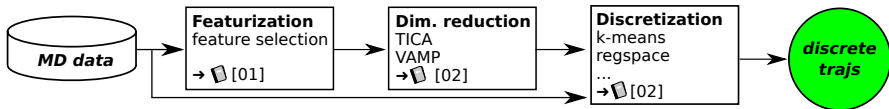


Figure: Workflow: state space discretisation



- ▶ All readers are Python-“iterable”, which means you can process data in chunks. The more general concept in PyEMMA is called ‘DataSource’.

```
1 | my_source = pyemma.coordinates.source(['traj001.xtc', ...])
2 | for element in my_source:
3 |     print(element)
```

Supported reader data formats:

- ▶ MD-simulation data (XTC, DCD, ... via MDTraj)
- ▶ NumPy (.npy) files
- ▶ Tabulated ASCII data (around three times more efficient than Numpy.loadtxt)
- ▶ Fragmented trajectories [(‘sim_0_part0.xtc’, ‘sim_0_part1.xtc’), ‘sim_1_part0.xtc’, ‘sim_1_part1.xtc’)]



Python package for reading/writing and analyzing molecular trajectories.
Analysis functions:

- ▶ distances
- ▶ bonds/angles/dihedrals
- ▶ hydrogen bonding identification
- ▶ secondary structure assignment
- ▶ NMR observables
- ▶ ... and many more

Supported formats:

- ▶ DCD
- ▶ XTC
- ▶ TRR
- ▶ PDB
- ▶ XYZ
- ▶ binpos
- ▶ NetCDF
- ▶ LH5
- ▶ HDF5
- ▶ ...

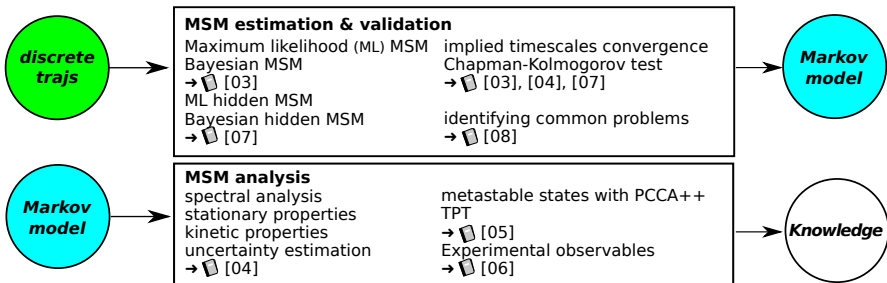


Figure: MSM estimation and analysis workflow.



Step	Goal	API function <small>(all in pyemma.msm package)</small>
1.a	choose lag time	<code>its = timescales_msm(dtrajs)</code>
1.b	choose lag time (visual inspection)	<code>pyemma.plots.plot_implied_timescales(its)</code>
2	estimate a model	<code>msm_obj = estimate_markov_model(dtrajs, lag)</code>
3.a	validate model	<code>ck_obj = msm_obj.cktest()</code>
3.b	validate model (vis. inspection)	<code>pyemma.plots.plot_cktest(ck_obj)</code>
4.a	Analyze slow processes	<code>msm_obj.timescales()</code> etc.
4.b	Perform coarse graining	<code>coarsed = msm_obj.pcca()</code>
4.c	Transition path analysis	<code>coarsed.tpt()</code>



Software overview and design patterns

- Python

- Anaconda stack

- Package overview

 - Coordinates package

 - MSM package

PyEMMA Development

- Principles

- Processes

 - GitHub

 - Continous Integration Services

- Collaboration



- ▶ Use Python as the glue to faster languages (C/C++, Fortran)
- ▶ Stable and easy to use high level user interface
- ▶ Open source (GNU Lesser Public license 3+, minimal restrictions on redistribution)
- ▶ Open development process on GitHub (everybody can contribute)
- ▶ Focus on speed and stability (NumPy, SciPy under the hood)
- ▶ Focus on good documentation (see <http://emma-project.org>)





- ▶ GitHub as frontend (collect issues/bugs, discuss proposed changes, plan new features, ...)
- ▶ Continuous integration/deployment (Travis-CI, AppVeyor, custom Jenkins instances)
- ▶ Unit-tests for API and implementation
- ▶ Integration tests of notebooks
- ▶ Release bug fixes regularly
- ▶ Release major/minor versions, if API changes.
- ▶ Preserve API compatibility (deprecate functions first, to notice users, that in the future their program/scripts will not work the same way as before)



- ▶ Before a release we freeze acceptance of new features (their milestone gets postponed to the next release)
- ▶ Testing sessions - eliminate all found bugs
- ▶ Deploy source archive to PyPI (installable with pip) and binaries to Anaconda.org binary services.
- ▶ Version scheme:
Major.minor.micro

major = major new (and API break features)

minor = new features preserving existing API

micro = patches/bug fixes



This repository [Pull requests](#) [Issues](#) [Gist](#) 🔔 + 👤

[markovmodel / PyEMMA](#) 👁 Unwatch 25 ★ Unstar 65 🍴 Fork 41

[Code](#) [Issues 76](#) [Pull requests 7](#) [Projects 0](#) [Wiki](#) [Pulse](#) [Graphs](#) [Settings](#)

Python API for Emma's Markov Model Algorithms <http://pyemma.org> [Edit](#)

[python](#) [hidden-markov-model](#) [molecular-dynamics](#) [analysis](#) [markov-state-model](#) [tica](#) [time-series](#) [Manage topics](#)

[4,552 commits](#) [7 branches](#) [32 releases](#) [20 contributors](#) [LGPL-3.0](#)

Branch: [devel](#) [New pull request](#) [Create new file](#) [Upload files](#) [Find file](#) [Clone or download](#)

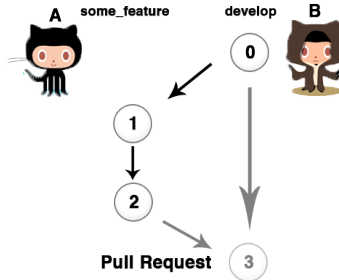
marscher committed on GitHub Merge pull request #1043 from marscher/prog_rep_no_class ... Latest commit df584ca 4 days ago
.github [github] added templates for new issues and pull requests 4 months ago
devtools Merge branch 'devel' into fix_plot_network a month ago
doc [doc] amend changelog 4 days ago
pyemma-ipython @ cdb59e9 [msm-bpti] do not import matplotlib explicitly a month ago
pyemma pytest cleanup 4 days ago
.gitattributes [versioneer] updated to v0.13 2 years ago

Figure: PyEMMA GitHub page



Collaboration on GitHub

1. Propose a change/feature via an issue
2. Create a local branch in Git to work on
3. Push the (tested) branch to your fork
4. Open a “pull request” (PR) on main repository (markovmodel/PyEMMA)
5. Discuss changes, eventually add more commits
6. Maintainer merges your PR



Propose file change on GitHub



...continued





- ▶ Create a GitHub account to directly post issues (**preferred**).
- ▶ Join our channel on Gitter.im
- ▶ Send mails to the developers (more overhead for us, might not reach somebody in time).



Thank you for your attention!
Further questions?