

# 將軍

Sören Sonnenburg<sup>†</sup>, Gunnar Rätsch<sup>‡</sup>, Fabio De Bona<sup>‡</sup>



**Fraunhofer** Institut  
Rechnerarchitektur  
und Softwaretechnik



MAX-PLANCK-GESELLSCHAFT

- <sup>†</sup> Fraunhofer FIRST.IDA, Kekuléstr. 7, 12489 Berlin, Germany
- <sup>‡</sup> Friedrich Miescher Laboratory of the Max Planck Society,  
Spemannstr. 39, 72076 Tübingen, Germany

## Machine Learning Toolbox SHOGUN features algorithms:

- to learn 2-class classification and regression problems
- to train hidden markov models
- toolbox's focus is on kernel methods esp. Support Vector Machines (SVMs) for computational biology
- also implements a number of linear methods like Linear Discriminant Analysis (LDA), Linear Programming Machine (LPM), (Kernel) Perceptrons

## FEATURES

- SHOGUN interfaces to Matlab<sup>TM</sup>, Octave and Python and R
- provides generic SVM object interfacing to several different SVM implementations, among them the state-of-the-art LibSVM and SVM<sup>light</sup>
- SVMs can be trained using a variety of common kernels (Linear, Polynomial, Gaussian) and recent String Kernels (TOP, Fisher, Locality Improved, Spectrum, Weighted DegreeKernel (with shifts))
- kernels can be combined; weighting can be learned using Multiple Kernel Learning.
- input feature-objects can be dense, sparse or strings and of type int/short/double/char; can be converted into different feature types.
- multiprocessor parallelization  $\Rightarrow$  able train on **10 million** examples

... and many more...

## GENERIC DEMO:

- Support Vector Classification
  - Task: separate 2 clouds of gaussian distributed points in 2D
  - octave, R, python
- Support Vector Regression
  - Task: learn a sine function
  - octave, R, python
- Hidden Markov Model
  - Task: 3 loaded dice are drawn 1000 times, find out when which dice was drawn
  - octave, R, python

## BIOINFORMATICS DEMO:

- Position Independent (e.g. Tissue Classification using Promotor Region)

```
AAACAAAA CGTAACTAATCTTTTAGAGAGAACGTTTCAACCATTTTGAG
AAGATTA ACTCATCACAGATTTT CATTACATACAGATATAATTCAA AATT
CACTCCCAAATCAACGATATTTAAAA TCACTAACACATCCGTCTGTGC
```

- Task: separate DNA strings, '-' class random ACGT, '+' class contains 'AAAAA' motif

- Position Dependent (e.g. Splice Site Classification)

```
AAACAAATAAGTAACTAATCTTTTAA GAAGAACGTTTCAACCATTTTGAG
AAGATTA AAAAAAAAAACAAATTTT AACATTACAGATATAATAATCTAATT
CACTCCCAAATCAACGATATTTTAA TTCACTAACACATCCGTCTGTGCC
```

- Task: separate DNA strings, '-' class random ACGT, '+' class 'AA' in the middle

- Mixture Position Dependent/Independent (e.g. Promoter Classification)

```
AAACAAATAAGTAACTAATCTTTTAA AGAGAACGTTTCAACCATTTTGAG
AAGATTA AAAAAAAAAACAAATTTTCAA TAAATACAGATATAATAATCTAATT
CACTCCCAAATCAACGATATTTAAA TTCACTAACACATCCGTCTGTGC
```

- Task: separate DNA strings, '-' class random 'ACGT', '+' class 'AAA' in the middle shifted  $\pm 15$

## DRAWBACKS AND FUTURE

### Now:

- Interface to R, octave, matlab, python using same/very similar syntax

```
sg("set_features", "TRAIN", traindat)          sg('set_features', 'TRAIN', traindat)
sg("set_labels", "TRAIN", trainlab)           sg('set_labels', 'TRAIN', trainlab)
sg("send_command", "set_kernel GAUSSIAN REAL 40 1") sg('send_command', 'set_kernel GAUSSIAN REAL 40 1')
sg("send_command", "init_kernel TRAIN")       sg('send_command', 'init_kernel TRAIN')
sg("send_command", "new_svm LIGHT")          sg('send_command', 'new_svm LIGHT')
sg("send_command", "c 10.0")                 sg('send_command', 'c 10.0')
sg("send_command", "svm_train")              sg('send_command', 'svm_train')
sg("set_features", "TEST", testdat)          sg('set_features', 'TEST', testdat)
sg("send_command", "init_kernel TEST")       sg('send_command', 'init_kernel TEST')
out=sg("svm_classify")                       out=sg('svm_classify')
```

- No “objects,” i.e. not possible with  $> 1$  SVM classifier

### Future:

- Focus on a new python object oriented interface (C++ classes are directly wrapped via swig)
- $\Rightarrow$  much cleaner: Demo

## SUMMARY

- SHOGUN is a large scale machine learning toolbox  
⇒ able to train on **10 million** examples
- unified SVM framework + many string kernels suitable for comp. biology
- Algorithms: HMM, LDA, LPM, Perceptron, SVM, SVR + many kernels
- Interfaces to matlab, octave, python, R

**Contribute + get support for all interfaces automagically!**

**Help us!**

- Documentation
- Testing
- Examples
- Test Suite

**Source Code is freely available under the GPLv2.**

<http://www.fml.tuebingen.mpg.de/raetsch/projects/shogun>