# Machine Learning Open Source Software and Benchmark Repository

Mikio L. Braun
TU Berlin
mikiobraun.de

October 7, 2010
Validation in Statistics and Machine Learning
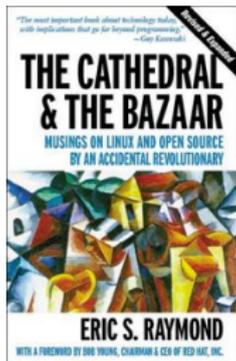WIAS Berlin

# Validation in Machine Learning and Statistics

| Machine Learning | Statistics |
|---|---|
| Solve hard computational tasks | Obtain scientific insight from data |

$\rightarrow$ Data and validation important both for ML and statistics, but for different reasons:

- **Machine Learning:** to share learning problems and compare existing methods.
- **Statistics:** to ensure that scientific insights are correct.

## "Open Source"



- ▶ Actually, open source is about a collaborative process to develop software (not unlike science!)
- ▶ Infrastructure: Source code revision system, bug trackers, mailing lists, discussion forums... .
- ▶ Once you release your code, you should enter this process.
- ▶ Opportunity for much faster interaction with "users"!

# Legal Implications of Reproducible Research

Victoria Stodden, *"Reproducible Research in Computational Science: Problems and Solutions For Data and Code Sharing"* http://videolectures.net/icml2010_stodden_rric/

- ▶ Releasing Source Code is not the same as Open Source!
- ▶ By default, "original expressions of ideas" is copyrighted (protects reproduction and derivative works, limited lifetime)
- ▶ Data? "Raw facts" not copyrightable. "Original selection and arrangement" is. Best option is to release to the public domain.

# Open Source Licenses

- ▶ Main purpose: Allow derived work.
- ▶ Differences: Derived work must also be released as open source, commercial use allowed, patents allowed, etc.
- ▶ Applying such a license: As easy as downloading the code and adding it as a file called LICENSE, adding links to that code.

# More on Licenses

- ▶ "Classical" source code license:
    - ▶ GNU Public License
    - ▶ BSD license,
    - ▶ Apache 2.0
    - ▶ "Lesser" GPL,
    - ▶ Affero GPL,
    - ▶ see www.opensource.org/licenses/alphabetical
- ▶ Creative Commons (creativecommons.org):
    - ▶ CC BY (attribution)
    - ▶ CC NC (no commercial use)
    - ▶ CC ND (no derived works)
    - ▶ CC SA (derived work must use same license)
- ▶ Public Domain (CC0): Waive all rights

# The Reproducible Research Standard

Victoria Stodden, *"Enabling Reproducible Research: Open Licensing For Scientific Innovation"* International Journal of Communications Law and Policy, Issue 13, 2009.

- Remove copyright's barrier to reproducible research,
- Realign the IP framework with longstanding scientific norms. A suite of license recommendations for computational science:

  1. Release media components (text, figures) under CC BY,
  2. Release code components under Modified BSD or similar,
  3. Release data to public domain (CC0) or attach an attribution license.

# Machine Learning Open Source Software

- MLOSS: NIPS Workshop 2006
- Position paper: S. Sonnenburg, M. L. Braun, C. S. Ong, S. Bengio, L. Bottou, G. Holmes, Y. LeCun, K.-R. Mller, F. Pereira, C. E. Rasmussen, G. Rtsch, B. Schlkopf, A. Smola, P. Vincent, J. Weston, R. Williamson, *The Need for Open Source Software in Machine Learning*, Journal of Machine Learning Research, 8(Oct):2443-2466, 2007
- mloss.org: Machine Learning Open Source Software
- JMLR MLOSS Track
- mldata.org: Machine Learning Data Set Repository

# mloss.org



- ▶ Open directory of machine learning related open source software projects.
- ▶ 272 projects registered so far.
- ▶ About 250 visitors per day.

Soeren Sonnenburg, Cheng Soon Ong, Mikio Braun

# JMLR MLOSS Special Track

http://jmlr.csail.mit.edu/mloss/



- ▶ Submit software together with 4 page description.
- ▶ About 20 projects published so far.

Soeren Sonnenburg, Cheng Soon Ong, Mikio Braun

# mldata.org



- Repository of data sets.
- Fully versioned, editable like a Wiki.
- Defines own data format based on HDF5. If used, additional features are available (e.g. automatic evaluation of prediction errors, download of data in other formats)