

# MACHINE LEARNING QUICK REFERENCE: BEST PRACTICES

Topic	Common Challenges	Suggested Best Practice
<b>Data Preparation</b>		
<b>Data collection</b>	<ul style="list-style-type: none"> <li>Biased data</li> <li>Incomplete data</li> <li>The curse of dimensionality</li> <li>Sparsity</li> </ul>	<ul style="list-style-type: none"> <li>Take time to understand the business problem and its context</li> <li>Enrich the data</li> <li>Dimension-reduction techniques</li> <li>Change representation of data (e.g. COO)</li> </ul>
<b>“Untidy” data</b>	<ul style="list-style-type: none"> <li>Value ranges as columns</li> <li>Multiple variables in the same column</li> <li>Variables in both rows and columns</li> </ul>	Restructure the data to be “tidy” by using the melt and cast process
<b>Outliers</b>	<ul style="list-style-type: none"> <li>Out-of-range numeric values and unknown categorical values in score data</li> <li>Undue influence on squared loss functions (e.g. regression, GBM, and <i>k</i>-means)</li> </ul>	<ul style="list-style-type: none"> <li>Robust methods (e.g. Huber loss function)</li> <li>Discretization (binning)</li> <li>Winsorizing</li> </ul>
<b>Sparse target variables</b>	<ul style="list-style-type: none"> <li>Low primary event occurrence rate</li> <li>Overwhelming preponderance of zero or missing values in target</li> </ul>	<ul style="list-style-type: none"> <li>Proportional oversampling</li> <li>Inverse prior probabilities</li> <li>Mixture models</li> </ul>
<b>Variables of disparate magnitudes</b>	<ul style="list-style-type: none"> <li>Misleading variable importance</li> <li>Distance measure imbalance</li> <li>Gradient dominance</li> </ul>	Standardization
<b>High-cardinality variables</b>	<ul style="list-style-type: none"> <li>Overfitting</li> <li>Unknown categorical values in holdout data</li> </ul>	<ul style="list-style-type: none"> <li>Discretization (binning)</li> <li>Weight of evidence</li> <li>Leave-one-out event rate</li> </ul>
<b>Missing data</b>	<ul style="list-style-type: none"> <li>Information loss</li> <li>Bias</li> </ul>	<ul style="list-style-type: none"> <li>Discretization (binning)</li> <li>Imputation</li> <li>Tree-based modeling techniques</li> </ul>
<b>Strong multicollinearity</b>	Unstable parameter estimates	<ul style="list-style-type: none"> <li>Regularization</li> <li>Dimension reduction</li> </ul>
<b>Training</b>		
<b>Overfitting</b>	High-variance and low-bias models that fail to generalize well	<ul style="list-style-type: none"> <li>Regularization</li> <li>Noise injection</li> <li>Partitioning or cross validation</li> </ul>
<b>Hyperparameter tuning</b>	Combinatorial explosion of hyper-parameters in conventional algorithms (e.g. deep neural networks, Super Learners)	<ul style="list-style-type: none"> <li>Local search optimization, including genetic algorithms</li> <li>Grid search, random search</li> </ul>
<b>Ensemble models</b>	<ul style="list-style-type: none"> <li>Single models that fail to provide adequate accuracy</li> <li>High-variance and low-bias models that fail to generalize well</li> </ul>	<ul style="list-style-type: none"> <li>Established ensemble methods (e.g. bagging, boosting, stacking)</li> <li>Custom or manual combinations of predictions</li> </ul>
<b>Model Interpretation</b>	Large number of parameters, rules, or other complexity obscures model interpretation	<ul style="list-style-type: none"> <li>Variable selection by regularization (e.g. L1)</li> <li>Surrogate models</li> <li>Partial dependency plots, variable importance measures</li> </ul>
<b>Computational resource exploitation</b>	<ul style="list-style-type: none"> <li>Single-threaded algorithm implementations</li> <li>Heavy reliance on interpreted languages</li> </ul>	<ul style="list-style-type: none"> <li>Train many single-threaded models in parallel</li> <li>Hardware acceleration (e.g. SSD, GPU)</li> <li>Low-level, native libraries</li> <li>Distributed computing, when appropriate</li> </ul>
<b>Deployment</b>		
<b>Model deployment</b>	Trained model logic must be transferred from a development environment to an operational computing system to assist in organizational decision making processes	<ul style="list-style-type: none"> <li>Portable scoring code or scoring executables</li> <li>In-database scoring</li> <li>Web service scoring</li> </ul>
<b>Model decay</b>	<ul style="list-style-type: none"> <li>Business problem or market conditions have changed since the model was created</li> <li>New observations fall outside domain of training data</li> </ul>	<ul style="list-style-type: none"> <li>Monitor models for decreasing accuracy</li> <li>Update/retrain models regularly</li> <li>Champion-challenger tests</li> <li>Online updates</li> </ul>

# MACHINE LEARNING QUICK REFERENCE: RESOURCES

## Publications

Statistical Modeling, The Two Cultures – Leo Breiman

- <http://projecteuclid.org/euclid.ss/1009213726>

Fifty Years of Data Science – David Donoho

- <http://courses.csail.mit.edu/18.337/2015/docs/50YearsDataScience.pdf>

Pattern Recognition and Machine Learning – Christopher Bishop

- <https://www.cs.princeton.edu/courses/archive/spring07/cos424/papers/bishop-regression.pdf>

Machine Learning with SAS Enterprise Miner – SAS White Paper

- [http://www.sas.com/content/dam/SAS/en\\_us/doc/whitepaper1/machine-learning-with-sas-enterprise-miner-107521.pdf](http://www.sas.com/content/dam/SAS/en_us/doc/whitepaper1/machine-learning-with-sas-enterprise-miner-107521.pdf)

An Overview of Machine Learning with SAS® Enterprise Miner™ - 2014 SGF Paper (SAS313-2014)

- <http://support.sas.com/resources/papers/proceedings14/SAS313-2014.pdf>

## Posts

An Introduction to Machine Learning – Patrick Hall on sas.com

- <http://blogs.sas.com/content/sascom/2015/08/11/an-introduction-to-machine-learning/>

7 Common Mistakes of Machine Learning – Cheng-Tao Chu on KDNuggets

- <http://www.kdnuggets.com/2015/03/machine-learning-data-science-common-mistakes.html>

How to build a deep neural network in SAS Enterprise Miner – Answer on SAS Data Mining community

- <https://communities.sas.com/t5/SAS-Communities-Library/How-to-build-a-deep-learning-model-in-SAS-Enterprise-Miner/ta-p/231190>

## Repos

A curated list of awesome Machine Learning frameworks, libraries and software

- [github.com/josephmisiti/awesome-machine-learning](https://github.com/josephmisiti/awesome-machine-learning)

Benchmark tests/results for open source implementations of the top machine learning algorithms

- [github.com/szilard/benchm-ml](https://github.com/szilard/benchm-ml)

Code/materials for integrating SAS with popular open source analytics technologies like Python and R.

- [github.com/sassoftware/enlighten-integration](https://github.com/sassoftware/enlighten-integration)

Quick reference tables for machine learning best practices and algorithm usage

- [github.com/sassoftware/enlighten-apply/tree/master/ML\\_tables](https://github.com/sassoftware/enlighten-apply/tree/master/ML_tables)

Library of SAS Enterprise Miner process flow diagrams to help you learn by example

- [github.com/sassoftware/dm-flow](https://github.com/sassoftware/dm-flow)