# MACHINE LEARNING QUICK REFERENCE: BEST PRACTICES

| Topic | Common Challenges | Suggested Best Practice |
|---|---|---|
| **Data Preparation** | | |
| **Data collection** | • Biased data<br>• Incomplete data<br>• The curse of dimensionality<br>• Sparsity | • Take time to understand the business problem and its context<br>• Enrich the data<br>• Dimension-reduction techniques<br>• Change representation of data (e.g., COO) |
| **"Untidy" data** | • Value ranges as columns<br>• Multiple variables in the same column<br>• Variables in both rows and columns | Restructure the data to be "tidy" by using the melt and cast process |
| **Outliers** | • Out-of-range numeric values and unknown categorical values in score data<br>• Undue influence on squared loss functions (e.g., regression, GBM, $k$-means) | • Robust methods (e.g., Huber loss function)<br>• Discretization (binning)<br>• Winsorizing |
| **Sparse target variables** | • Low primary event occurrence rate<br>• Overwhelming preponderance of zero or missing values in target | • Proportional oversampling<br>• Inverse prior probabilities<br>• Mixture models |
| **Variables of disparate magnitudes** | • Misleading variable importance<br>• Distance measure imbalance<br>• Gradient dominance | Standardization |
| **High-cardinality variables** | • Overfitting<br>• Unknown categorical values in holdout data | • Discretization (binning)<br>• Weight of evidence<br>• Leave-one-out event rate |
| **Missing data** | • Information loss<br>• Bias | • Discretization (binning)<br>• Imputation<br>• Tree-based modeling techniques |
| **Strong multicollinearity** | Unstable parameter estimates | • Regularization<br>• Dimension reduction |
| **Training** | | |
| **Overfitting** | High-variance and low-bias models that fail to generalize well | • Regularization<br>• Noise injection<br>• Partitioning or cross validation |
| **Hyperparameter tuning** | Combinatorial explosion of hyperparameters in conventional algorithms (e.g., deep neural networks, super learners) | • Local search optimization, including genetic algorithms<br>• Grid search, random search |
| **Ensemble models** | • Single models that fail to provide adequate accuracy<br>• High-variance and low-bias models that fail to generalize well | • Established ensemble methods (e.g., bagging, boosting, stacking)<br>• Custom or manual combinations of predictions |
| **Model Interpretation** | Large number of parameters, rules, or other complexity obscures model interpretation | • Variable selection by regularization (e.g., L1)<br>• Surrogate models<br>• Partial dependency plots, variable importance measures |
| **Computational resource exploitation** | • Single-threaded algorithm implementations<br>• Heavy reliance on interpreted languages | • Train many single-threaded models in parallel<br>• Hardware acceleration (e.g. SSD, GPU)<br>• Low-level, native libraries<br>• Distributed computing, when appropriate |
| **Deployment** | | |
| **Model deployment** | Trained model logic must be transferred from a development environment to an operational computing system to assist in organizational decision-making processes | • Portable scoring code or scoring executables<br>• In-database scoring<br>• Web service scoring |
| **Model decay** | • Business problem or market conditions have changed since the model was created<br>• New observations fall outside domain of training data | • Monitor models for decreasing accuracy<br>• Update/retrain models regularly<br>• Champion-challenger tests<br>• Online updates |

https://github.com/sassoftware/enlighten-apply/tree/master/ML_tables

§sas