

MACHINE LEARNING QUICK REFERENCE: BEST PRACTICES

Topic	Common Challenges	Suggested Best Practice
Data Preparation		
Data collection	<ul style="list-style-type: none">Biased dataIncomplete dataThe curse of dimensionalitySparsity	<ul style="list-style-type: none">Take time to understand the business problem and its contextEnrich the dataDimension-reduction techniquesChange representation of data (e.g. COO)
“Untidy” data	<ul style="list-style-type: none">Value ranges as columnsMultiple variables in the same columnVariables in both rows and columns	Restructure the data to be “tidy” by using the melt and cast process
Outliers	<ul style="list-style-type: none">Out-of-range numeric values and unknown categorical values in score dataUndue influence on squared loss functions (e.g. regression, GBM, and <i>k</i>-means)	<ul style="list-style-type: none">Robust methods (e.g. Huber loss function)BinningWinsorizing
Sparse target variables	<ul style="list-style-type: none">Low primary event occurrence rateOverwhelming preponderance of zero or missing values in target	<ul style="list-style-type: none">Proportional oversamplingInverse prior probabilitiesMixture models
Variables of disparate magnitudes	<ul style="list-style-type: none">Misleading variable importanceDistance measure imbalanceGradient dominance	Standardization
High-cardinality variables	<ul style="list-style-type: none">OverfittingUnknown categorical values in holdout data	<ul style="list-style-type: none">BinningWeight of evidenceLeave-one-out event rate
Missing data	<ul style="list-style-type: none">Information lossBias	<ul style="list-style-type: none">BinningImputationTree-based modeling techniques
Strong multicollinearity	Unstable parameter estimates	<ul style="list-style-type: none">RegularizationDimension reduction
Training		
Overfitting	High-variance and low-bias models that fail to generalize well	<ul style="list-style-type: none">RegularizationNoise injectionPartitioning or cross validation
Hyperparameter tuning	Combinatorial explosion of hyper-parameters in conventional algorithms (e.g. deep neural networks, Super Learners)	<ul style="list-style-type: none">Local search optimization, including genetic algorithmsGrid search, random search
Ensemble models	<ul style="list-style-type: none">Single models that fail to provide adequate accuracyHigh-variance and low-bias models that fail to generalize well	<ul style="list-style-type: none">Established ensemble methods (e.g. bagging, boosting, stacking)Custom or manual combinations of predictions
Model Interpretation	Large number of parameters, rules, or other complexity obscures model interpretation	<ul style="list-style-type: none">Variable selection by regularization (e.g. L1)Surrogate modelsPartial dependency plots, variable importance measures
Computational resource exploitation	<ul style="list-style-type: none">Single-threaded algorithm implementationsHeavy reliance on interpreted languages	<ul style="list-style-type: none">Train many single-threaded models in parallelHardware acceleration (e.g. SSD, GPU)Low-level, native librariesDistributed 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	<ul style="list-style-type: none">Portable scoring code or scoring executablesIn-database scoringWeb service scoring
Model decay	<ul style="list-style-type: none">Business problem or market conditions have changed since the model was createdNew observations fall outside domain of training data	<ul style="list-style-type: none">Monitor models for decreasing accuracyUpdate/retrain models regularlyChampion-challenger testsOnline updates

