# Machine Learning Modelling in R : : CHEAT SHEET

## Supervised & Unsupervised Learning

| | ALGORITHM | DESCRIPTION | R PACKAGE::FUNCTION | SAMPLE CODE |
|---|---|---|---|---|
| **SUPERVISED LEARNING** | **NBC** Naïve Bayes classifier | A classification technique based on Bayes' Theorem with an assumption of independence among predictors. In simple terms, a Naïve Bayes classifier assumes that the presence of a particular feature in a class is unrelated to the presence of any other feature | e1071::naiveBayes | naiveBayes(class ~ ., data = x) |
| | **kNN** k–Nearest Neighbours | A non-parametric method used for classification and regression. In both cases, the input consists of the k closest training examples in the feature space. The output depends on whether k-NN is used for classification or regression | class::knn | knn(train, test, cl, k = 1, l = 0, prob = FALSE, use.all = TRUE) |
| | **LREG** Linear Regression | Model the linear relationship between a scalar dependant variable Y and one or more explanatory variables (or independent variables) denoted X | stats::lm | lm(dist ~ speed, data=cars) |
| | **LREG** Logistic Regression | Used to predict a binary outcome (1 / 0, Yes / No, True / False) given a set of independent variables. | stats::glm | glm(Y ~ ., family = binomial (link = 'logit'), data = X) |
| | **TM** Tree-Based Models | The idea is to consecutively divide (branch) the training dataset based on the input features until an assignment criterion with respect to the target variable into a "data bucket" (leaf) is reached | rpart::rpart | rpart(Kyphosis ~ Age + Number + Start, data = kyphosis) |
| | **ANN** Artificial Neural Network | Neural networks are built from units called perceptrons. Perceptrons have one or more inputs, an activation function and an output. An ANN model is built up by combining perceptrons in structured layers. | neuralnet::neuralnet | neuralnet(f,data=train_,hidden=c(5,3),linear.output=T) |
| | **SVM** Support Vector Machine | A data classification method that separates data using hyperplanes | e1071::svm | svm(formula, data = NULL, ..., subset, na.action = na.omit, scale = TRUE) |
| **UNSUPERVISED LEARNING** | **PCA** Principal Component Analysis | A procedure that uses an orthogonal transformation to convert a set of observations of possibly correlated variables into a set of values of linearly uncorrelated variables called principal components. | stats::prcomp stats::princomp FactoMineR::PCA ade4::dudi.pca amap::acp | **stats** : prcomp(formula, data = NULL, subset, na.action, ...) **stats** : princomp(formula, data = NULL, subset, na.action, ...) **FactoMineR** : PCA(decathlon, quanti.sup = 11:12, quali.sup=13) **ade4** : dudi.pca(deug$tab, center = deug$cent, scale = FALSE, scan = FALSE) **amap** : acp(lubisch) |
| | **kMC** k-Mean Clustering | Aims at partitioning n observations into k clusters in which each observation belongs to the cluster with the nearest mean | stats::kmeans | kmeans(x, centers, iter.max = 10, nstart = 1, algorithm = c("Hartigan-Wong", "Lloyd", "Forgy", "MacQueen"), trace=FALSE) |
| | **HCL** Hierarchical Clustering | An approach which builds a hierarchy from the bottom-up, and doesn't require the number of clusters to be specified beforehand. | stats::hclust | hclust(d, method = "complete", members = NULL) |

## Meta-Algorithm, Time Series & Model Validation

| | ALGORITHM | DESCRIPTION | R PACKAGE::FUNCTION | SAMPLE CODE |
|---|---|---|---|---|
| **META ALGORITHM** | **REGU** Regularisation L1 (Lasso) L2 (Ridge) | Regularization adds a penalty on the different parameters of a model to reduce the freedom of the model. Hence, the model will be less likely to fit the noise of the training data and will improve the generalization abilities of the model | glmnet::glmnet | **L1** : glmnet(myMatrixA, myMatrixB , family = "gaussian", alpha = 1) **L2** : glmnet(myMatrixA, myMatrixB , family = "gaussian", alpha = 0) |
| | **BOO** Boosting | A process of iteratively refining, e.g. by reweighting, of estimated regression and classification functions (though it has primarily been applied to the latter), in order to improve predictive ability. | **Paramatric model** - mboost::glmboost | glmboost(Yen ~ ., data = curr1[trnidxs,]) |
| | **BAG** Bagging | Bagging is a way to increase the power of a predictive statistical model by taking multiple random samples (with replacement) of the training data set, and using each of them to construct a separate model and separate predictions for the original test set | **All models**: foreach **Tree models**: ipred::bagging | **foreach** : d <- data.frame(x=1:10, y=rnorm(10)) s <- foreach(d=iter(d, by='row'), .combine=rbind) %dopar%d identical(s, d) **ipred**: bagging(formula, data, subset, na.action=na.rpart, \dots) |
| | **PRU** Pruning | Pruning is a technique that reduces the size of decision tree by removing sections of the tree that provide little power to classify instances. Pruning reduces the complexity of the final classifier and hence improves predictive accuracy by reducing overfitting | rpart::prune | prune(x, cp = 0.1) |
| | **RFO** Random Forrest | An ensemble learning method for classification, regression and other tasks, that operate by constructing a multitude of decision trees at training time and outputting the class that is the mode of the classes (classification) or mean prediction (regression) | randomForest::randomForest | randomForest(X ~ ., data = Y , subset = mySub) |
| **TIME SERIES** | **STS** Lead-lag analysis, Auto-correlation, Spectral analysis, Time series clustering, Seasonality, Trend.... | Random sampling of observations for training and testing a model can be an issue when faced with a times dimension. Random sampling may either destroy serial correlation properties in the data which we would like to exploit | stats xts forecast spectral TTR .... | **Auto-correlation**: acf(x, lag.max = NULL, type = c("correlation", "covariance", "partial")) **Spectral Analysis**: spec.pgram(myTs, spans = NULL) **Seasonal Decomposition of Time Series** - stl(x, s.window = 7, t.window = 50, t.jump = 1) ...... |
| **MODEL VALIDATION** | **PM** Performance metrics | Depends on the problem: • **Regression**: squared errors, outliers, error rate... • **Classification**: Accuracy, precision, recall, F-score... | **Regression**-stats::outlierTest, stats::qqPlot ... **Classification**-ROCR:: **Tree**: caret:: confusionMatrix | **Regression**: fit = lm(Y~X,data=myData) outlierTest(fit) qqPlot(fit, main="QQ Plot") |
| | **BVT** Biais-Variance Tradeoff | • Simple models with few parameters are easier to compute but may lead to poorer fits (**high bias**) • Complex models may provide more accurate fits but may over-fit the data (**high variance**) | Tailored to the analysis | Tailored to the analysis |
| | **CV** Cross validation | Cross validation compares the test performances of different model realisations with different sets or values of parameters | caret::createDataPartition caret::createFolds | createDataPartition(classes, p = 0.8, list = FALSE) |
| | **LC** Learning Curves | Learning curves plot a model's training and test errors, or the chosen performance metric, depending on the training set size | caret:: learing_curve_dat | learing_curve_dat(dat, outcome = NULL, proportion = (1:10)/10, test_prop = 0, verbose = TRUE, ...) |

## Standard Modelling Workflow



## Time Series View



TRAINING VALIDATION PRODUCTION

**A ① ** A simple model is trained using a first sub sample of the original data set

**A ② ** A meta model is evaluated using a first sub sample of the original data set

**① ** Simple models can be aggregated to define a meta models.

**② ** Simple models are validated using a second sub sample of the original dataset

**③ ** Meta models are validated using a second sub sample of the original dataset

**④ ** Validated models are tested in production

**⑤ ** Depending results, simple and meta models might go through the training and validation steps again until production results are acceptable