

Demo Steps: Using the R API to Create, Score, and Assess Models

From Jupyter Lab, select **File Browser > Home > Courses > EVMLOPRC > Notebooks** and select the **R_Machine_Learning_Demo.ipynb** notebook. Scroll to Part 2.

1. Partition the **HMEQ** data set into 70% for training and 30% for validation.

```
cas.sampling.srs(conn,
  table = indata,
  sampct = 70,
  seed = 919,
  partind = TRUE,
  output = list(casOut = list(name = indata, replace = T),
    copyVars = 'ALL')
)
```

Selected arguments:

Argument	Description
table	specifies the in-memory input data table. The table is used in combination with the CAS session connection object.
sampct	specifies the sample percentage to be used for sampling or partitioning.
seed	specifies the integer to use to start the pseudorandom number generator.
partind	when set to <i>true</i> , generates a partition indicator column in the output table.
output	creates, on the server, a table that contains the sample output or partition output.
casOut	specifies the settings for an output table.
name	specifies the in-memory table name to save the action results.
replace	when set to <i>true</i> , replaces the table with the results of the action.
copyVars	specifies a list of one or more variables to be copied from the input table to the output table. The keyword ALL copies all variables.

The action added only a single binary partition indicator variable to the data set by using the `partind=TRUE` argument (`partind = 1` for training and `partind = 0` for validation).

2. Use SQL to find the count and percent of each level of the new **_PartInd_** variable.

```
loadActionSet(conn, 'fedSql')
```

```
counts = cas.fedSql.execDirect(conn, query =
  "
  SELECT _PartInd_, count(*)
  FROM hmeq
  GROUP BY _PartInd_;
  "
)$`Result Set`
```

counts

The **counts** table is copied and brought to the client. Native open source syntax can act on the table further. In general, SQL can be used equivalently to PROC SQL. Under the `fedSql` action set and `execDirect` action, SQL takes advantage of the distributed environment and parallel processing.

```
counts$Percent = counts$COUNT/sum(counts$COUNT)
```

counts

The result from the query is a data frame. Open source code is then used to add the **percent** column to the table. The percent is equivalent to the sampling proportion from the `srs` action.

3. Using the **defCasTable** function, refresh the object reference to the **HMEQ** data set and then use the **SWAT** function mean to find the proportion of training cases specified by the **_PartInd_** variable.

```
castbl = defCasTable(conn, table=indata)
mean(castbl$`_PartInd_`)
```

4. Use the super action from the **varReduce** action set to perform supervised dimension reduction. The action identifies a set of variables that jointly explain the maximum amount of variance that is contained in the response variables. Find the number of variables that are needed to explain 90% of the variability in the response.

```
loadActionSet(conn, 'varReduce')
```

```
varReduce_obj = cas.varReduce.super(conn,
  table = indata,
  target = target,
  inputs = inputs,
  nominals = nominals,
  varexp = 0.90
)
```

names(varReduce_obj)

Selected arguments:

Argument	Description
table	specifies the in-memory input data table. The table is used in combination with the CAS session connection object.
target	specifies the target variable to use for analysis.
inputs	specifies the input variables to use for analysis.

nominals	specifies the nominal variables to use for analysis.
varexp	specifies the fraction of the total variance to be explained.

```
varReduce_obj$SelectionSummary
```

```
varReduce_obj$SelectedEffects
varReduce_obj$SelectedEffects$Variable
```

The super action selected all 12 inputs to explain 90% of the variability in the target. For larger data sets, it can be helpful to find a subset of inputs before beginning the modeling phase of the analysis.

The super action selects variables by using the discriminant criterion that is specified in linear discrimination analysis.

5. Use the logistic action from the regression action set to train a logistic regression model.

```
loadActionSet(conn, 'regression')
```

```
cas.regression.logistic(conn,
  table = list(name = indata, where = '_PartInd_ = 1'),
  classVars = nominals[-1],
  model = list(depvar=target, effects=inputs, dist='binomial',
    link='logit'),
  store = list(name='lr_model', replace=TRUE)
)
```

Selected arguments:

Argument	Description
table	specifies the in-memory input data table. The table is used in combination with the CAS session connection object.
where	specifies a conditional argument for the observations to be used in the analysis.
classVars	names the classification variables for the analysis.
model	names the dependent variable, explanatory variable, and model options.
store	saves the regression model.

Note: The logistic action can also perform regularization and stepwise selection.

6. Use the svmTrain action from the svm action set to train a support vector machine.

```
loadActionSet(conn, 'svm')
```

```
cas.svm.svmTrain(conn,
  table = list(name = indata, where = '_PartInd_ = 1'),
  target = target,
  inputs = inputs,
  nominals = nominals,
  kernel = 'polynomial',
  degree = 2,
  savestate = list(name = 'svm_model', replace = TRUE)
)
```

Selected arguments:

Argument	Description
kernel	specifies the kernel type.
Degree	specifies the degree of the polynomial kernel.
savestate	specifies the table to save the model for future scoring.

7. Load the decisionTree action set. Notice that it provides all the actions to train and score all three tree-based models. Use the dtreeTrain action to train a decision tree.

```
loadActionSet(conn, 'decisionTree')
```

```
cas.decisionTree.dtreeTrain(conn,
  table = list(name = indata, where = '_PartInd_ = 1'),
  target = target,
  inputs = inputs,
  nominals = nominals,
  casOut = list(name = 'dt_model', replace = TRUE)
)
```

Selected arguments:

Argument	Description
table	specifies the in-memory input data table.
target	specifies the target variable to use for analysis.
inputs	specifies the input variables to use for analysis.
nominals	specifies the nominal variables to use for analysis.
casOut	specifies the table to store the decision tree model. If it is not specified, a random name is generated.

The model information displays the default parameters for the decision tree model. These can be changed by using arguments within the action.

8. Use the forestTrain action to train a random forest and use 1000 trees in the forest.

```
cas.decisionTree.forestTrain(conn,
  table = list(name = indata, where = '_PartInd_ = 1'),
  target = target,
  inputs = inputs,
```

```

    nominals = nominals,
    nTree     = 1000,
    casOut    = list(name = 'rf_model', replace = TRUE)
)

```

Selected argument:

Argument	Description
nTree	specifies the number of trees to create.

9. Use the `gbtreeTrain` action to train a gradient-boosting model and use 1000 tree iterations.

```

cas.decisionTree.gbtreeTrain(conn,
  table = list(name = indata, where = '_PartInd_ = 1'),
  target = target,
  inputs = inputs,
  nominals = nominals,
  nTree = 1000,
  casOut = list(name = 'gbt_model', replace = TRUE)
)

```

10. Load the `neuralNet` action set and use the `annTrain` action to train a neural network. Use a single hidden layer with 150 neurons.

```
loadActionSet(conn, 'neuralNet')
```

```

cas.neuralNet.annTrain(conn,
  table = list(name = indata, where = '_PartInd_ = 1'),
  target = target,
  inputs = inputs,
  nominals = nominals,
  hidden = list(150),
  nloOpts = list(optm1Opt = list(maxIters = 100,
                                fConv = 1e-10),
                lbfgsOpt = list(numCorrections = 6)),
  casOut = list(name = 'nn_model', replace = TRUE)
)

```

Selected arguments:

Argument	Description
hidden	specifies the number of neurons for each hidden layer. For example, <code>hidden={100,50}</code> specifies two hidden layers, one with 100 neurons and the other with 50.
nloOpts	specifies the nonlinear optimization options.
maxIters	specifies the maximum iterations allowed for optimization. The default is 10.
fConv	specifies a stopping value when the objective functions fail to change more than a value in the allotted iterations. The default is 1e-05.
lbfgsOpt	specifies options for the Broyden-Fletcher-Goldfarb-shanno algorithm, an iterative method for solving unconstrained nonlinear optimization problems.
numCorrections	specifies the number of corrections used in the LBFGS update. The default is 20.
casOut	specifies the in-memory table to store the model.

11. Using the previously saved models, score each model on the validation data. The support vector machine is saved as an analytical store and the other models are saved as data tables. Therefore, the `svm` model requires the `aStore` action set to score the model, and the other models have their own scoring actions.

```
loadActionSet(conn, 'aStore')
```

```

#Score the support vector machine model
cas.aStore.score(conn,
  table = list(name = indata, where = '_PartInd_ = 0'),
  rstore = "svm_model",
  out = list(name="svm_scored", replace=TRUE)
)

```

Selected arguments:

Argument	Description
rstore	specifies the input analytical store.
out	specifies the in-memory table to store the scored data.

```

#Score the logistic regression model
lr_score_obj = cas.regression.logisticScore(conn,
  table = list(name = indata, where = '_PartInd_ = 0'),
  restore = "lr_model",
  casout = list(name="lr_scored", replace=TRUE),
  copyVars = list(target)
)

```

```

#Score the decision tree model
dt_score_obj = cas.decisionTree.dtreeScore(conn,
  table = list(name = indata, where = '_PartInd_ = 0'),
  model = "dt_model",
  casout = list(name="dt_scored", replace=TRUE),
  copyVars = list(target),
  encodename = TRUE,
  assessorerow = TRUE
)

```

```
#Score the random forest model
```

```

rf_score_obj = cas.decisionTree.forestScore(conn,
  table = list(name = indata, where = '_PartInd_ = 0'),
  model = "rf_model",
  casout = list(name="rf_scored",replace=TRUE),
  copyVars = list(target),
  encodename = TRUE,
  assessonerow = TRUE
)

#Score the gradient boosting model
gb_score_obj = cas.decisionTree.gbtreeScore(conn,
  table = list(name = indata, where = '_PartInd_ = 0'),
  model = "gbt_model",
  casout = list(name="gbt_scored",replace=TRUE),
  copyVars = list(target),
  encodename = TRUE,
  assessonerow = TRUE
)

#Score the neural network model
nn_score_obj = cas.neuralNet.annScore(conn,
  table = list(name = indata, where = '_PartInd_ = 0'),
  model = "nn_model",
  casout = list(name="nn_scored",replace=TRUE),
  copyVars = list(target),
  encodename = TRUE,
  assessonerow = TRUE
)

```

Selected arguments:

Argument	Description
model	specifies the table that contains the model information.
copyVars	specifies the variables to transfer from the input table to the output table.
encodename	specifies whether encoding the variable names in the generated output table adds the target variable name to the predicted probability variable name, such as the predicted probabilities of each response variable level.
assessonerow	when set to <i>true</i> , causes predicted probabilities to be added to the results table for the event levels. All event probabilities are included as separate columns and are used with the assess action.

```

nn_score_obj$OutputCasTables
nn_score_obj$ScoreInfo

```

12. From the percentile action set, use the assess action to evaluate the scoring results for each model.

```
loadActionSet(conn, 'percentile')
```

```

# Change the name of the prediction variable for logistic regression
cas.dataStep.runCode(conn, code="
  data lr_scored;
  set lr_scored;
  rename _PRED_ = P_BAD1;
run;
"
)

```

```

# Add the target variable to the svm scored table
cas.dataStep.runCode(conn, code="
  data svm_scored;
  merge svm_scored(keep=P_BAD1) lr_scored(keep=BAD);
run;
"
)

```

The first DATA step simply changes the variable name for the predicted probabilities from `_PRED_` to `P_BAD1` to be consistent with the naming convention of the other models. The default variable name is the prefix `P_` followed by the target and level, **`P_BAD1`**.

```

# Create prediction variable name
assess_input = paste("P_", target, "1", sep = "")

```

```

# Assess the logistic regression model
lr_assess_obj = cas.percentile.assess(conn,
  table = 'lr_scored',
  inputs = assess_input,
  casout = list(name="lr_assess",replace=TRUE),
  response = target,
  event = "1"
)

```

```

# Assess the support vector machine model
svm_assess_obj = cas.percentile.assess(conn,
  table = 'svm_scored',
  inputs = assess_input,
  casout = list(name="svm_assess",replace=TRUE),
  response = target,
  event = "1"
)

```

```

# Assess the decision tree model
dt_assess_obj = cas.percentile.assess(conn,
  table = "dt_scored",
  inputs = assess_input,

```

```

    casout = list(name="dt_assess",replace=TRUE),
    response = target,
    event = "1"
)

# Assess the random forest model
rf_assess_obj = cas.percentile.assess(conn,
  table = "rf_scored",
  inputs = assess_input,
  casout = list(name="rf_assess",replace=TRUE),
  response = target,
  event = "1"
)

#Assess the gradient boosting model
gb_assess_obj = cas.percentile.assess(conn,
  table = "gbt_scored",
  inputs = assess_input,
  casout = list(name="gbt_assess",replace=TRUE),
  response = target,
  event = "1"
)

# Assess the neural network model
nn_assess_obj = cas.percentile.assess(conn,
  table = "nn_scored",
  inputs = assess_input,
  casout = list(name="nn_assess",replace=TRUE),
  response = target,
  event = "1"
)

```

Selected arguments:

Argument	Description
event	specifies the formatted value of the response variable that represents the event.
inputs	specifies the input variable for the analysis. This is the predicted probability of the modeling level for binary classification.
response	specifies the original response variable for the model.
casOut	specifies the in-memory table to save the results.

```
nn_assess_obj$outputCasTables
```

Notice that the assess action saves two separate tables. The first is the name that is specified by the casOut argument and the second is the same name but it ends in _ROC. Each table saves different information to assess the model.

```
cas.table.fetch(conn, table='nn_assess', to=5)
```

```
cas.table.fetch(conn, table='nn_assess ROC', to=5)
```

13. Use the defCasTable function to create an object reference to the ROC tables for each model. Then use the to.casDataFrame function to download the tables to analyze the results on the client.

```

lr_assess_ROC = defCasTable(conn, tablename = "lr_assess_ROC")
lr_assess_ROC = to.casDataFrame(lr_assess_ROC)
lr_assess_ROC$Model = 'Logistic Regression'

svm_assess_ROC = defCasTable(conn, tablename = "svm_assess_ROC")
svm_assess_ROC = to.casDataFrame(svm_assess_ROC)
svm_assess_ROC$Model = 'Support Vector Machine'

dt_assess_ROC = defCasTable(conn, tablename = "dt_assess_ROC")
dt_assess_ROC = to.casDataFrame(dt_assess_ROC)
dt_assess_ROC$Model = 'Decision Tree'

rf_assess_ROC = defCasTable(conn, tablename = "rf_assess_ROC")
rf_assess_ROC = to.casDataFrame(rf_assess_ROC)
rf_assess_ROC$Model = 'Random Forest'

gbt_assess_ROC = defCasTable(conn, tablename = "gbt_assess_ROC")
gbt_assess_ROC = to.casDataFrame(gbt_assess_ROC)
gbt_assess_ROC$Model = 'Gradient Boosting'

nn_assess_ROC = defCasTable(conn, tablename = "nn_assess_ROC")
nn_assess_ROC = to.casDataFrame(nn_assess_ROC)
nn_assess_ROC$Model = 'Neural Network'

```

A new variable, **Model**, was added to each local data frame to specify the model name for the results.

14. Stack the ROC data frames into one data frame and then print the confusion matrix results at a 0.50 cutoff rate for the probability of an event.

```

df_assess = rbind(lr_assess_ROC, svm_assess_ROC, dt_assess_ROC, rf_assess_ROC, gbt_assess_ROC, nn_assess_ROC)
cutoff_index = df_assess[,3]==0.5
compare = df_assess[cutoff_index,]
rownames(compare) = NULL
compare[,c('Model','TP','FP','FN','TN')]

```

15. Use the **_ACC_** variable to print and compute the misclassification for each model.

```

compare$Misclassification = 1 - compare$'_ACC_'
miss = compare[order(compare$Misclassification),
  c('Model','Misclassification')]
rownames(miss) = NULL
miss

```

16. Using ggplot, compare ROC curves and add the area under the curve to the plot legend.

```

# Add a new column to be used as the ROC curve label
df_assess$Models = paste(df_assess$Model,
                          round(df_assess$'_C_', 3), sep = ' - ')

#Subset the data frame with only three variables
df_roc = df_assess[c('FPR', '_Sensitivity_', 'Models')]
colnames(df_roc) = c("FPR", "Sensitivity", "Models")

# Create the ROC curve
ggplot(data = df_roc,
       aes(x = FPR, y = Sensitivity, colour = Models)) +
  geom_line() +
  labs(x = 'False Positive Rate', y = 'True Positive Rate')

```

17. Create table references for the other assess results and download them to the client.

```

lr_assess_lift = defCasTable(conn, tablename = "lr_assess")
lr_assess_lift = to.casDataFrame(lr_assess_lift)
lr_assess_lift$Model = 'Logistic Regression'

svm_assess_lift = defCasTable(conn, tablename = "svm_assess")
svm_assess_lift = to.casDataFrame(svm_assess_lift)
svm_assess_lift$Model = 'Support Vector Machine'

dt_assess_lift = defCasTable(conn, tablename = "dt_assess")
dt_assess_lift = to.casDataFrame(dt_assess_lift)
dt_assess_lift$Model = 'Decision Tree'

rf_assess_lift = defCasTable(conn, tablename = "rf_assess")
rf_assess_lift = to.casDataFrame(rf_assess_lift)
rf_assess_lift$Model = 'Random Forest'

gbt_assess_lift = defCasTable(conn, tablename = "gbt_assess")
gbt_assess_lift = to.casDataFrame(gbt_assess_lift)
gbt_assess_lift$Model = 'Gradient Boosting'

nn_assess_lift = defCasTable(conn, tablename = "nn_assess")
nn_assess_lift = to.casDataFrame(nn_assess_lift)
nn_assess_lift$Model = 'Neural Network'

df_assess = rbind(lr_assess_lift, svm_assess_lift, dt_assess_lift, rf_assess_lift, gbt_assess_lift, nn_assess_lift)

```

18. Using ggplot, create a lift curve.

```

df_lift = df_assess[c('_Depth_', '_CumLift_', 'Model')]
colnames(df_lift) = c("Depth", "CumLift", "Models")

ggplot(data = df_lift,
       aes(x = Depth, y = CumLift, colour = Models)) +
  geom_line() +
  labs(x = 'Population Percentage', y = 'Lift')

```

19. Sample 75% of the HMEQ data set and then bring the data to the client. For larger data sets, you need to reduce the amount of data by creating a sample subset before you bring it to the client so that you do not use too much RAM. Then fit a gradient boosting model locally.

```

cas.sampling.srs(conn,
  table = indata,
  samp_pct = 75,
  seed = 12345,
  partind = FALSE,
  output = list(casOut = list(name = 'mysam', replace = T),
                copyVars = 'ALL')
)

# Bring data locally
mysam = defCasTable(conn, table='mysam')
df = to.casDataFrame(mysam)
df = df[,c(target, inputs, '_PartInd_')]

# Create dummy variables through one-hot encoding
df.dum = df[,nominals[-1]]
dummies = dummyVars('~ .', data = df.dum)
df.ohe = as.data.frame(predict(dummies, newdata = df))
df.all.combined = cbind(df[, -c(which(colnames(df) %in%
  nominals[-1]))], df.ohe)

# Split into training and validation
train = df.all.combined[df.all.combined['_PartInd_'] == 1,]
valid = df.all.combined[df.all.combined['_PartInd_'] == 0,]

# Train the XGBoost model
set.seed(101112)
bst = xgboost(
  data = data.matrix(train[, -1]),
  label = data.matrix(train[, 1]),
  nround = 25,
  objective = "binary:logistic",
  eta = 0.1,
  max_depth = 5,
  subsample = 0.5,
  colsample_bytree = 0.5
)

```

20. Print the misclassification rate for the local model and then combine and print misclassification rates for all models, both local and CAS.

```

pred = as.numeric(predict(bst, data.matrix(valid[,-1]),
                        missing = 'NAN') > 0.5)
Misclassification = mean(as.numeric(pred > 0.5) != valid[,1])
xgb = data.frame(cbind(Model = 'R - XGBoost', Misclassification))
xgb

err = data.frame(rbind(miss, xgb))
err[, -1] = round(as.numeric(as.character(err[, -1])), 7)
err = err[order(err[, -1]), ]
rownames(err) = NULL
err

```

Because the CAS actions results are local, they can be combined with the local modeling results to create a data frame for comparison.

21. Use the `tableInfo` action from the `table` action set to view all the tables that are currently on the server.

```
cas.table.tableInfo(conn)$TableInfo[, 1:3]
```

Note: Unless they are saved or promoted, all tables are deleted after the CAS session is terminated.

22. Use the `caslibInfo` action to view the path for the default caslib **CASUSER**. Then use the `addCaslib` action to create a user-defined caslib that is named **mycl** located in the home directory.

```
cas.table.caslibInfo(conn, active=FALSE, caslib="casuser")
```

```

cas.table.addCaslib(conn, name="mycl", path=Sys.getenv("HOME"),
  dataSource="PATH", description="Personal File Save Location",
  activeOnAdd = FALSE)

```

Selected arguments:

Argument	Description
name	specifies the name of the caslib to be added.
dataSource	specifies the data source type.
path	specifies the data source specific information.
description	specifies a string description of the new caslib.
activeOnAdd	When it is set to <i>true</i> , the new caslib becomes the active caslib. When it is set to <i>false</i> , the default caslib remains active.

The new caslib is **local**. This means that the connection to the directory is terminated when the session ends. However, the data that is saved in the library is permanent. The new caslib is not active, which means that the name must be specified so that it can be used. Otherwise, the default caslib is assumed.

```
cas.table.caslibInfo(conn)
```

23. Use the `save` action to save the gradient boosting model on the server in the **mycl** caslib. Also, use the `attribute` action to save the table attributes for the model.

```

cas.table.save(conn, caslib = 'mycl', table = list(name =
  'gbt_model'), name = 'best_model_gbt', replace = T)

cas.table.attribute(conn, caslib = 'CASUSER', table =
  'gbt_model_attr', name = 'gbt_model', task='convert')

```

```

cas.table.save(conn, caslib = 'mycl', table = 'gbt_model_attr',
  name = 'attr', replace = T)

```

Selected argument:

Argument	Description
task	specifies the task to perform. The convert task creates a table for the attributes so that they can be saved.

24. Using the `dropCaslib` action, drop the user-defined caslib, **mycl**.

```
cas.table.dropCaslib(conn, caslib="mycl")
```

25. Promote the **HMEQ** data set to the global scope. This causes the table to persist on the server unless it is explicitly dropped. It can be accessed from other APIs and other sessions.

```

cas.table.promote(conn, caslib="casuser", name=indata)
cas.table.tableInfo(conn)

```

26. End the session and disconnect from CAS.

```
cas.session.endSession(conn)
```