

LVI Linear Discriminant Analysis

December 18, 2012

Prepared for:

Forest Analysis and Inventory Branch

BC Ministry of Forests, Lands and Natural Resource Operations

By

Ian Moss, RPF, PhD

Tesera Systems Inc.

Victoria, BC

LVI Linear Discriminant Analysis

Table of Contents

1. Introduction	1
2. Create an R working directory	1
3. Start with Variable Selection from a Large Number of Variables	2
4. Run Discriminant Analysis for a Single X-Variable Set	4
5. Run Linear Discriminant Analysis for Multiple X-Variable Sets.....	6
6. Produce Unique X-Variable Subset Correlation Matrix	9
7. Produce Original Classification Unique X- or Y-Variable Subset Box and Scatter Plots.....	10
8. Combine Evaluation Datasets	11
9. Alternative Variable Selection Procedure (Under Development).....	12

LVI Linear Discriminant Analysis

1. Introduction

Assumes windows machine

All code tested using R 2.9.1

Assumes competency with computers and some low level familiarity with Python and R

R packages installed (see appendix)

All operating within an R working directory (Rwd)

Scripts

Data processing

Statistical analyses

Data visualization

Data interpretation (summary statistics)

Modular construction for ease of extension

Don't need to program in R

Can help if you want to learn how to program in R – some working code with documentation

2. Create an R working directory

1. Create an R working directory on your computer as follows: \\Rwd.
2. Copy the Rwd contents into the same directory you just established in the previous step; the Rwd directory should have three subdirectories: Rdata, Rdocs, and RScript.
3. Start up a new R session.
4. click **File, Open Script** at the top RHS of the R interpreter and navigate to the Rwd directory.
5. Navigate to the //Rwd// RScript file and open the SetRwd.R script file.
6. Change setwd("E:\\Rwd") to indicate the address of Rwd on your computer, e.g. setwd("C:\\Rwd"); note that you must use "\\" as subfile separators in the file address.
7. Click **Edit, Run all**.
8. A new file will now be stored in your Rwd directory called OpenSession.RData.
9. If you now close the R session, go to the Rwd directory, double click on OpenSession.RData a new R session will open up in the interpreter.
10. Now if you once again click on **File, Open** you will see that the interpreter automatically opens up in the Rwd directory.

3. Start with Variable Selection from a Large Number of Variables

1. Load Dataset

1.1. LoadDatasetAndAttachVariableNames.R (DATA is LVINEW.csv)

2. Select and declare (numerical) classification (Y-) variable as a “Factor”

2.1. DeclareClassificationVariableAsFactor.R

Note that this is an opportunity to change the classification Y-variable dataset

3. Select (independent) X variable subset

3.1. SelecXVariableSubset_v1 (DATA is XVARSELV1.csv)

3.2. SelectXVariableSubset_v2 (DATA is XVARSELV2.csv)

Note that the second version (v2) is designed to accommodate multiple variable selection; however, when used in this procedure only header and the first row may be entered, otherwise only the last variable set in the list will be used in the next step. The first version (XVARSELV1.txt) was originally designed to identify the X variable superset from which specific variable sets would be extracted. However this could also be accomplished using the second version (XVARSEV2.txt) format.

(The first version is used currently as the starting point – it was also constructed manually and was intended for initial variable selection)

4. Load R library package, subselect.

4.1. Loadsubselect-R-Package.R

(Note if you wish to view the data and make changes to it you can run the following script:

ViewLviNewDataset.R

This will open up a new screen and allow you to view the data. However, if this is deployed using a large dataset it may cause the interpreter to become non-responsive since it tries to write all of the data into the interpreter once the data editor is closed.

5. Run IdaHmat in subselect package

5.1. RunLinearDiscriminantAnalysis_subselect_IdaHmat.R

6. Select chosen criteria and run IdaHmat variable selection routine

6.1. Run-IdaHmat-VariableSelection-Improve.R

LVI Linear Discriminant Analysis

Note that for variable selection you can select a range in the number of variables to be used in the model by changing the minimum (minNvar) and the maximum (maxNvar). You can also control the number of solutions that you would like to investigate for each number of variables (between minNvar and maxNvar, inclusive). Finally you can set different criteria for variable selection as follows:

- Roy's first root statistic ("ccr12")
- Wilks' Lamda ("Wilkes")
- Chi squared ("x12")
- And the Zeta 2 coefficient ("zeta2")

Currently of the alternatives are listed in the script file. To deselect a choice put a # in front of the line (R recognizes these as comments). To select a criterion, remove the # sign from the front of the line. Note that if two of the criteria are selected, the last one listed (toward the bottom of the script) will be the one used.

7. Extract all variable subsets derived from step 6 and put in data frame called SOLSUM (solution summary).

7.1. ExtractVariableNameSubsets.R

8. Write SOLSUM from step 7 to VARSELECT.csv in R working directory

8.1. WriteDataframeToCsvFile.R

9. Python Code used here.

9.1. Using a Python module create a reformatted list of all of the unique combinations of variables and print it to a file called XVARSELV in the LVI directory; Run the following routine:

9.1.1. EXTRACT_RVARIABLE_COMBOS.py

OUTPUT (".csv" comma delimited files):

XVARSELV contains the list of unique combinations of variables developed from running the variable selection routine lda.Hmat.

UNIQUEVAR contains a unique list of variable names compiled from all variable sets produced in lda.Hmat.

Note that the output is assigned to the following directory:

"E\\Rwd\\"

4. Run Discriminant Analysis for a Single X-Variable Set

1. Load Dataset

1.1. LoadDatasetAndAttachVariableNames.R (DATA is LVINEW.txt)

2. Select and declare (numerical) classification variable as a "Factor"

2.1. DeclareClassificationVariableAsFactor.R

Note that this is an opportunity to change the classification Y-variable if there are a number of different classifications you would like to investigate.

3. Set prior classification distribution as being uniform or as per sample

3.1. Uniform: ComputeUniformPriorClassificationProbabilityDistribution.R

3.2. Sample: ComputeSamplePriorClassProbabilityDistribution.R

4. Select (independent) X variable subset

4.1. SelectXVariableSubset_v1 (DATA is XVARSELV1.csv)

4.2. SelectXVariableSubset_v2 (DATA is XVARSELV2.csv)

These are loaded in the //Rwd//

Note that the second version (v2) is designed to accommodate multiple variable selection; however, when used in this procedure only header and the first row may be entered, otherwise only the last variable set in the list will be used in the next step.

5. Load Discriminant Analysis R-package

5.1. LoadMASS-R-Package.R

6. Run Linear Discriminant Analysis

6.1. RunLinearDiscriminantAnalysis_MASS_Ida.R

6.2. RunLinearDiscriminantAnalysis_MASS_Ida_TakeOneLeaveOne.R

Notes:

In *RunLinearDiscriminantAnalysis_MASS_Ida.R* the Take-One-Leave-One option is disabled (CV = FALSE). As a result the following output is available following the discriminant analysis (note by typing the command, indicated in bold, the results can be printed out in the interpreter:

LVI Linear Discriminant Analysis

- lvi.Ida\$prior** this produces the prior probability distribution (established in step 3 above) used to represent the distribution of observations amongst the classes (CLASSIFICATION).
- lvi.Ida\$counts** the number of observations by class.
- lvi.Ida\$means** the mean for each X-variable by class
- lvi.Ida\$scaling** the discriminant functions are scaled so that the mean z-score for each function is 0. Note that this is equivalent to subtracting observed X-variable value from the mean for each of the variables and then multiplying by the discriminant functions.
- lvi.Ida\$svd** the ratio's of between to within-group standard deviations in the linear discriminant variables. These are also referred to as eigenvalues. The squares of these figures are the canonical F-statistics. When the squares of these figures are converted into proportions of the total – this is equivalent to the proportion of the total (Between-to-within) variance explained by each discriminant function.
- lvi.Ida\$N** is the number of observations contained in the dataset.

Within this same routine the following additional output is also available:

- class.pred\$class** (the predict function) produces the class assignments to each observation (with all observations used in the discriminant analysis); also at the bottom of this output, the unique class names (or numbers) are listed.
- class.pred\$posterior** the estimated posterior probability distributions based on the prior distribution calculated as follows (see Hora and Wilcox 1982, Dillon and Goldstein 1984 pp. 392 – 393).
- class.pred\$x** the scores for each if the test cases associated with each variate (function)
- class.table** the (table function produces a) classification contingency table with the original class distribution in rows and the predicted class distribution in columns.

In *RunLinearDiscriminantAnalysis_MASS_Ida_TakeOneLeaveOne.R* the Take-One-Leave-One option is enabled (CV =TRUE). The following output may be obtained:

- lvi.Ida\$class** the class assigned to each observation
- lvi.Ida\$posterior** the posterior probabilities developed using Take-One-Leave-One; these are superior to those not involving the Take-One-Leave-One process (Bates and Wilcox 1982; Dillon and Goldstein 1984, pp. 406-409).
- class.table** the (table function produces a) classification contingency table with the original class distribution in rows and the predicted class distribution in columns.

5. Run Linear Discriminant Analysis for Multiple X-Variable Sets

1. Load Dataset

1.1. LoadDatasetAndAttachVariableNames.R (DATA is LVINEW.csv)

2. Select and declare (numerical) classification variable as a "Factor"

2.1. DeclareClassificationVariableAsFactor.R

Note that this is an opportunity to change the classification Y-variable if there are a number of different classifications you would like to investigate.

3. Compute prior classification distribution as being uniform or as per sample

3.1. Uniform: ComputeUniformPriorClassificationProbabilityDistribution.R

3.2. Sample: ComputeSamplePriorClassProbabilityDistribution.R

4. Write selected prior distribution to file

4.1. WritePriorDistributionToFile.R

OUTPUT (.csv file)

PRIOR Contains a list of classes (CLASS) and associated prior probabilities (PROIRD).

5. Load Discriminant Analysis R-package

5.1. LoadMASS-R-Package.R

6. Select (independent) X variable subset

6.1. SelectXVariableSubset_v2.1 (DATA is XVARSELV.csv)

The VARSELV2.txt is created by following procedure number 2 above, entitled, "Start with Variable Selection from a Large Number of Variables."

7. Run (Multiple) Linear Discriminant Analysis for Multiple Sets of Variables – Take-One-Leave-One

7.1. RunMultipleLinearDiscriminantAnalysis_MASS_Ida_TakeOneLeaveOne.R

Note that this routine uses the Take-One-Leave-One routine for the purpose of rating the quality of the variable sets in terms of their accuracies in classification based on producing contingency tables as 1 output, and in terms of the posterior estimation of error as another output.

LVI Linear Discriminant Analysis

Note that the posterior error of estimation is calculated following the procedures of Hora and Wilcox (1982; Equation 9):

$$\hat{e} = 1 - N^{-1} \sum_{i=1}^N \max [P(Y | X_i)] \quad \text{Eq. 1}$$

Where,

\hat{e} is the estimated (posterior) error

N Is the total number of observations

$P(Y | X_i)$ is the probability of class Y, where Y is equal to 1 to m classes, given a set of variables, X_i , where i equals 1 to n observations.

8. Write (Multiple Linear) Results from Step 6 to Files

8.1. WriteMultipleLinearDiscriminantAnalysis_MASS_Ida_TOLO_to_File.R

OUTPUT (.csv files)

CTABULATION: contains the contingency table data from which the Cohen's (1960) Coefficient of Agreement can be calculated.

POSTERROR: contains the errors estimated using Eq. 1 above.

9. Run (Multiple) Linear Discriminant Analysis for Multiple Sets of Variables – All Observations

9.1. RunMultipleLinearDiscriminantAnalysis_Mass_Ida.R

10. Write Multiple Linear results from step 8 to files.

10.1. WriteMultipleLinearDiscriminantAnalysis_MASS_Ida.R

OUTPUT (.csv files)

CTABALL: contains the contingency table data from which the Cohen's (1960) Coefficient of Agreement can be calculated.

VARMEANS: contains the mean values for each variable by class.

DFUNCT: contains the discriminant functions for each axis and combination of variables.

BWRATIO: contains the between-to-within variance ratio of the differences in class Z-statistics associated with each discriminant function (i.e. eigenvalues).

11. Compile Take-One-Leave-One CTABULATION Classification Accuracy Statistics

11.1. COHENS_KHAT.py

(Data is CTABULATION.csv)

LVI Linear Discriminant Analysis

Note that the input file was produced using step 7 in these procedures.

OUTPUT (.csv files in Rwd directory)

CTABSUM provides statistics as indicated in Table 2.

Table 2. A description of variables included in the output file: CTABSUM.

Variable	Name
VARSET	Variable Set
OA	Overall Accuracy
KHAT	Coefficient of Agreement
MINPA	Minimum Producer Accuracy
MAXPA	Maximum Producer Accuracy
MINUA	Minimum User Accuracy
MAXUA	Maximum User Accuracy

A variable set is associated with a given combination of different kinds and numbers of variables that as developed in steps 5, 6 and 7 above. The overall accuracy indicates the proportion of observations that were correctly classified according to the original (reference) classification. Cohen's coefficient of agreement is an indicator of the overall success rate after having removed the potential for a certain level of agreement to occur by chance.

The minimum producer accuracy indicates the minimum number of correctly classified observations given the total number of observations assigned to any given class, amongst all classes by way of discriminant analysis in this case, and expressed as a proportion. The maximum producer accuracy is similarly derived but with respect to the maximum. These figures are also related to the maximum and errors of omission amongst all of the classes (e.g. $1 - \text{MINPA}$, and $1 - \text{MAXPA}$).

The minimum user accuracy indicates the minimum number of correctly classified observations given the total number of observations as originally assigned to any given class, amongst all classes by way of discriminant analysis in this case, and expressed as a proportion. The maximum user accuracy is similarly derived but with respect to the maximum. These figures are also related to the maximum and errors of commission amongst all of the classes (e.g. $1 - \text{MINPA}$, and $1 - \text{MAXPA}$).

6. Produce Unique X-Variable Subset Correlation Matrix

Note that as a guideline you may wish to exclude any one variable in pairs with correlations > 0.7 (or < -0.7) a priori. The step of removing correlated variables may best be done

1. Load dataset

1.1. LoadDatasetAndAttachVariableNames.R (DATA is LVINEW.csv)

2. Select unique X variable subset

2.1. SelectUniqueXVariableSubset.R (DATA is UNIQUEVAR.csv)

3. Compile Unique Variable Correlation Matrix

3.1. CompileUniqueVariableCorrelationMatrixSubset.R

4. Create a unique variable correlation matrix file for printing

4.1. CreateUniqueVarCorrelationMatrixFileForPrinting.R

5. Select new X variable subset (v2.1)

5.1. SelectXVariableSubset_v2.1.R (DATA is XVARSELV21.csv)

6. Add variable subset indicators to correlation matrix file if option 5.1 selected

6.1. AddVariableSubsetIndicatorsToCorrelationMatrix.R

Note that this routine labels each variable indicator set as I1, I2, I3 ... in the order that they are produced in step 5.

7. Write unique variable correlation matrix to a file

7.1. WriteUniqueVarCorrelationMatrix.R

OUTUPT (.csv file)

UCORCOEF: This produces a table of correlation coefficients and indicator variables for each variable subset with 1's assigned to variable pairs that exist in the subset, and 0's assigned to all other variable pairs.

MINMAXCOR: This is a compilation of the maximum and minimum correlations within each variable set across all variable pairs (excluding identical pairs for

LVI Linear Discriminant Analysis

which the correlations are 1) within each variable subset. This is one criterion that can be used to select preferred sets of variables.

7. Produce Original Classification Unique X- or Y-Variable Subset Box and Scatter Plots

1. Load dataset

1.1. LoadDatasetAndAttachVariableNames.R (DATA is LVINEW.csv)

2. Declare (numerical) classification (Y-) variable as a “Factor”

2.1. DeclareClassificationVariableAsFactor.R
(Interaction of User)

3. Select unique X or Y variable subset

3.1. SelectUniqueXVariableSubset.R (DATA is UNIQUEVAR.csv)
3.2. SelectXVariableSubset_v1.R (DATA is XVARSELV1.csv)
3.3. SelectYVariableSubset_v1.R (DATA is XVARSELV1.csv)

4. Load R package – graphics

4.1. LoadGraphics-R-Package.R

5. Run the box plot script

If X variables have been selected then use:

5.1. CreateUniqueVariableClassificationBoxPlots.R

Else Y-variables were selected; use:

5.2. CreateYVariableClassificationBoxPlots.R

Note that you must right-click on each graph to proceed to the next graph. The first graph will be blank.

6. Run the scatter plot script

6.1. CreateUniqueVariableScatterPlots.R

Note that this can produce a lot of graphs. Specifically if there a total of K unique variables, then $(k*(K-1))/2$ graphs will be produced. For 20 variables that is equal to 190 graphs.

LVI Linear Discriminant Analysis

Having said that you can end the graphics session at any time, by closing the graphics box (right clicking on the “X” button in the top right hand corner of the window).

The figures illustrate the scatter of observations for y-axis versus x-axis variables. A regression line (red, y vs. x) and a lowess line (blue; a locally weighted polynomial regression line – similar to a moving average) are also shown on each figure to facilitate interpretation of the trends.

7. Produce X vs. Y variable scattergrams

If X variables have already been selected then select Y-variables:

7.1. SelectYVariableSubset_v1.R (DATA is XVARSELV1)

Else Y variables previously selected; Select X-variables:

7.2. SelectUniqueXVariableSubset.R (DATA is UNIQUEVAR.csv)

7.3. SelectXVariableSubset_v1.R (DATA is XVARSELV1.csv)

Then run the following script:

7.4. CreateYvXVariableScatterPlot.R

8. Combine Evaluation Datasets

Use a Python program, COMBINE_EVALUATION_DATASETS.py to combine the following datasets for purposes of overall assessment (OUTPUT: ASSESS.csv):

Table 3. A list of files combined into one file: ASSESS.csv.

Input	Description
PyRDataDict.csv	This contains the data types (e.g. string, float, integer) associated with each of the input tables, except XVARSELV
MINMAXCOR	This contains the minimum and maximum correlation coefficients amongst all pairs of variables contained within a variable set.
CTABSUM	See Table 2.
POSTERIOR	See Eq. 1 above.
XVARSELV	The actual variable sets produced during the variable selection process.

9. Alternative Variable Selection Procedure (Under Development)

1. Load dataset
 - 1.1. LoadDatasetAndAttachVariableNames.R (DATA is LVINEW.csv)
2. Declare (numerical) classification (Y-) variable as a “Factor”
 - 2.1. DeclareClassificationVariableAsFactor.R
3. Set prior classification distribution as being uniform or as per sample
 - 3.1. Uniform: ComputeUniformPriorClassificationProbabilityDistribution.R
 - 3.2. Sample: ComputeSamplePriorClassProbabilityDistribution.R
4. Select X variable subset (v1)
 - 4.1. SelectXVariableSubset_v1.R
5. Load klaR Package
 - 5.1. LoadklaR-R-Package.R
6. Load combinat Package
 - 6.1. LoadCombinat-R-Package.R
7. Run pairwise class variable selection
 - 7.1. RunLinearDiscriminantAnalysis_klaR_pvs.R

Note that this routine uses linear discriminant analysis (or alternatives such as quadratic or reduced discriminant procedures – making it more flexible than the subselect package which only provides for the standard linear discriminant analysis procedure). The procedure then compares each possible pair of classes in turn (using the pvs command in klaR) and selects the best variable sets according to certain criteria. The basic criteria are “stepclass” (forward, backward, or both), “ks.test” (Kolmogorov-Smirnov test), or “greedy.wilks” (Wilks’ lambda). This procedure could be applied within a bootstrap procedure to generate multiple variable sets for further testing.

LVI Linear Discriminant Analysis

Remaining on the list:

(Need to do)

1. Nearest Neighbour and RMSE Statistics (Both NN Classification with Centroids; pre and post Discriminant Classifications – and complete NN RMSE Statistics)
 - 1.1. Euclidean
 - 1.2. MSN (Mahalanobis)

(Like to do)

2. Classification of Target Dataset Using desired Prior and Posterior Distributions)
3. Apply Results to Target Dataset
4. Bootstrap Analyses
5. Use the concept of tolerance to eliminate collinearity – i.e. regress each X-variable on the remaining variables and remove those with a tolerance \leq a given threshold. Do so in order of least to most significant variables based on an ANOVA. This would be a useful screening tool to eliminate issues with Multicollinearity
6. Extend process to quadratic discriminant analysis and reduced discriminant analysis (use klaR package – a start has been made on this)
7. For graphical analysis – get unique pairs of variables used in equations – eliminate pairs that don't actually occur in any one equation (reduces number of graphs).
8. Determine optimal number of classes in classification given Y-variables (used to develop the classification and given X-variables.
9. Select Observation Subsets based on certain criteria.
10. Output actual versus predicted class for each observation for further graphical analysis

References

Cohen, J. 1960. A Coefficient of Agreement for Nominal Scales. *Educational and Psychological Measurement* XX(1): 37 – 46.

Dillon, W.R., and Goldstein, M. 1984. *Multivariate analysis*. Methods and applications. John Wiley & Sons, New York, NY, US.

Hora, S.C. and Wilcox, W.B. 1982. Estimation of error rates in several-population discriminant analysis. *Journal of Marketing Research* 19(1):57-61.

Additional References

Birkal, D. 2006. Regularized Discriminant Analysis. Lecture Notes. <http://www.uni-leipzig.de/~strimmer/lab/courses/ss06/seminar/slides/daniela-2x4.pdf> [accessed Dec 19 2012]

Davis, H.Z., Mesznik, R. and Lee, J.Y. 2009. Finding an internal optimum in the classification of management accounting information: The role of fuzzy sets. *Management Accounting* 17:203-216.

Liu, H. and Yu, L. 2005. Toward integrating feature selection algorithms for classification and clustering. *IEEE Transactions on Knowledge and Data Engineering* 17(4):491-502.

Appendix: A listing of Python Modules and R Scripts and R Working Directory Structure

Rwd Working Directory Structure

This is a dictionary created to maintain LVI data management processes. It is created separately from the location of the R interface, packages, etc.

1. E:
 - 1.1. Rwd
 - 1.1.2. Python
 - 1.1.2.1. Admin
 - 1.1.2.2. DATDICT
 - 1.1.2.2.1. PyRDataDict.csv
 - 1.1.2.3. PyReadError
 - 1.1.3. Rdata
 - 1.1.3.1. Archived
 - 1.1.4. Rdocs
 - 1.1.4.1. LOGISTIC
 - 1.1.4.2. MANOVA
 - 1.1.4.3. MDA
 - 1.1.4.4. Packages
 - 1.1.4.5. R-Language
 - 1.1.5. RScript
 - 1.1.5.1. Additional

R-packages

1. MASS
2. subselect
3. klaR
4. combinat

Note: to install a new packages you can use the following script:

installPackages.R

Enter the name of the package in quotation marks after the "packageName =" command and then run the script.

Python Module Library (\\Rwd\\Python\\)

1. COHENS_KHAT.py
2. COMBINE_EVALUATION_DATASETS.py
3. EXTRACT_RVARIABLE_COMBOS.py

LVI Linear Discriminant Analysis

Python Module Admin (Support) Files (\\Rwd\\Python\\Admin\\):

1. dictionaryDBUtilities.py
2. fileUtilities.py
3. innerJoinDictDBs.py
4. makeDD.py
5. PRINTv1.py
6. readCSV.py
7. typeDataset.py

Python Data Dictionary (\\Rwd\\Python\\DATDICT\\)

1. PyRDataDict.csv

DRAFT

LVI Linear Discriminant Analysis

R Script Library (\\Rwd\\RScript\\)

1. AddVariableSubsetIndicatorsToCorrelationMatrix.R
2. CompileUniqueVariableCorrelationMatrixSubset.R
3. ComputeSamplePriorClassProbabilityDistribution.R
4. ComputeUniformPriorClassificationProbabilityDistribution.R
5. CreateUniqueVarCorrelationMatrixFileForPrinting.R
6. CreateUniqueVariableClassificationBoxPlots.R
7. CreateUniqueVariableScatterPlots.R
8. CreateYVariableClassificationBoxPlots.R
9. CreateYvXVariableScatterPlots.R
10. DeclareClassificationVariableAsFactor.R
11. ExtractVariableNameSubsets.R
12. installPackages.R
13. LoadCombinat-R-Package.R
14. LoadDatasetAndAttachVariableNames.R
15. LoadGraphics-R-Package.R
16. LoadklaR-R-Package.R
17. LoadMASS-R-Package.R
18. Loadsubselect-R-Package.R
19. Run-IdaHmat-VariableSelection-Improve.R
20. RunLinearDiscriminantAnalysis_klaR_pvs.R
21. RunLinearDiscriminantAnalysis_MASS_Ida.R
22. RunLinearDiscriminantAnalysis_MASS_Ida_TakeOneLeaveOne.R
23. RunLinearDiscriminantAnalysis_subselect_IdaHmat.R
24. RunMultipleLinearDiscriminantAnalysis_Mass_Ida.R
25. RunMultipleLinearDiscriminantAnalysis_MASS_Ida_TakeOneLeaveOne.R
26. SelectUniqueXVariableSubset.R
27. SelectXVariableSubset_v2.1.R
28. SelectXVariableSubset_v2.R
29. SelectXVariableSubset_v1.R
30. SelectYVariableSubset_v1.R
31. SetRwd.R
32. ViewLviNewDataset.R
33. WriteDataframeToCsvFile.R
34. WriteMultipleLinearDiscriminantAnalysis_MASS_Ida_to_File.R
35. WriteMultipleLinearDiscriminantAnalysis_MASS_Ida_TOLO_to_File.R
36. WritePriorDistributionToFile.R
37. WriteUniqueVarCorrelationMatrix.R

LVI Linear Discriminant Analysis

The data (.csv) files used as input and output:

1. OpenSession.RData (after it has initially been created following procedure number 1 above)
2. Output (.csv) files in the Rwd directory:

ID	File	Source	Description
1	ASSESS	Python	File with summary statistics for assessment of results from modeling
2	BWRATIO	R	Ratios: Square root of between to within variance for each discriminant function
3	CTABALL	R	Cross tabulation of results using all data in model calibration.
4	CTABSUM	Python	Cross tabulation statistic summary derived from CTABULATION
5	CTABULATION	R	Cross tabulation results using Take-OneLeave-One
6	DFUNCT	R	Discriminant Functions
7	LVINEW	Manual	User defined base data including X, Y and CLASSIFICATION variables
8	LVINORM	Manual	User defined normalized (mean and standard deviation) data (continuous variables only) derived from LVINEW
9	MINMAXCORR	R	Minimum and Maximum Correlations within each X variable set
10	POSTERIOR	R	Error of estimation using posterior probabilities (see Hora and Wilcox 1982)
11	PRIOR	R	Prior probability selected by user
12	UCORCOEF	R	Correlation Coefficients amongst variables in UNIQUEVAR
13	UNIQUEVAR	Python	Unique variable set associated with variables listed in VARSELECT
14	VARMEANS	R	Variable means by original CLASS assignments
15	VARSELECT	R	String table describing variable sets used to derive XVARSELV
16	XVARSELV	Python	Variable sets used as further input into discriminant analysis
17	XVARSELV1	Manual	Initial file for selecting eligible X (and Y) variable datasets for use in analysis
18	XVARSELV2	Manual	Same format as XVARSELV but with only 1 variable set

Also note the PyDataDict.csv that is used in Python to designate variables as a certain type, e.g. integer, float, or string.