

Math 6365 first homework HW1

due date = Thursday feb 8th at midnight

This HW1 is a preliminary study of a concrete classification task of your choice, to prepare for “solving” this specific classification task in HW2 by implementation of MLPs

0)

select a publicly available data set associated to a clearly described CLASSIFICATION TASK,

1)

begin your report by writing your own description of the CLASSIFICATION TASK involved :

Outline a word description of “cases”; give the number N of cases

Give the number p of “descriptors” for each case. Outline the concrete meaning of each descriptor

Describe concretely each “class” $C_1 C_2 \dots C_k$ involved in the classification task. Give the size $|C_j|$ of each class C_j and the number k of classes. If k is larger than 5 , eliminate the classes $C_6 C_7 \dots C_k$ as well as all the cases corresponding to classes $C_6 C_7 \dots C_k$

Indicate how the correct classification of each case was originally obtained and how it is encoded in the data base

Indicate if there are any publication analyzing the classification task associated to this data set

2) A descriptor is called “discrete” or “categorical” if it can take only a “small” finite set of distinct values, which may be indexed by numerical values, or by distinct “names” in various formats. Of course one can always re-index this finite set of values by arbitrarily selected numbers. Examples : job type, accident severity, car type, etc.

A descriptor is called “continuous” if its set of potential values includes a full segment of real numbers. Examples: weight, length, width, volume, etc.

For each descriptor $D_1 \dots D_p$, indicate if it is categorical or continuous.

For each continuous descriptor compute and display its range , its min and max values, its mean and standard deviation. Compute and display its histogram

For each categorical descriptor indicate the finite list of distinct values , and the corresponding frequencies of these values in the whole data set

3)

For each descriptor D , and for each class C_j , compute the histogram $HIST_j$ of D values for the cases belonging to class C_j . If graphically feasible, display the k histograms $HIST_1 \dots HIST_k$ on a single graph. Interpret these informations to evaluate intuitively if the descriptor D will be helpful to discriminate between some pairs of classes.

One can use “Kolmogorov-Smirnov tests” to compare pairs of histograms HIST_i and HIST_j in order to decide if they are strongly distinct or not.

For each continuous descriptor D compute and display the mean M_j and standard deviation S_j of D values for the cases belonging to class C_j. One can then use classical “t-tests” to compare pairs of means M_i and M_j to decide if $|M_i - M_j|$ is “strongly” different from 0 or not ,

Both the t-tests and the Kolmogorov-Smirnov test are easy to use functions available in Matlab as well as in R

4)

For each case CAS_n with $n = 1, 2, \dots, N$ regroup all the continuous descriptors of CAS_n in a single vector V_n of fixed dimension r . Implement a PCA analysis for the cloud of N vectors V_1, V_2, \dots, V_N . Explain the PCA method in your report. Compute and display graphically the r eigenvalues $S_1 \geq S_2 \geq \dots \geq S_r$ generated by the PCA. Compute and display graphically the ratios $RAT_j = (S_1 + S_2 + \dots + S_j) / (S_1 + S_2 + \dots + S_N)$. Determine the smallest j such that $RAT_j \geq 0.95$.

5)

Call $W_1 \dots W_r$ the eigenvectors associated to S_1, S_2, \dots, S_r . Compute and display graphically the orthogonal projections of the N vectors V_1, V_2, \dots, V_N on the 3 dimensional vector space generated by the orthonormal basis $\{W_1, W_2, W_3\}$. Use different colors for these projected points namely 1 color per class. Interpret the results in your report

Do the same operation for the 3 separate projections onto the three planes $\{W_1, W_2\}$, $\{W_1, W_3\}$, and $\{W_2, W_3\}$. Interpret the results in your reports.