

KMLMM course. ZIP Practical work 2

Course: 2019-2020

Prof. Tomàs Aluja

We have normalized handwritten digits, automatically scanned from envelopes by the U.S. Postal Service in 16 x 16 grayscale images (from -1 to 1). Each line consists of the id (0-9) followed by the 256 grayscale values. We dispose of a training set of 7291 digits and a test set of 2007 digits. (files “zip_train.dat” and “zip_test.dat” respectively).

The purpose is to continue the exercise made session 1 using Multivariate Regression and a Principal Components Regression. Now we will try IBA as a component based methodology to predict the digits.

Steps for conducting the practice

1. Read the “zip_train.dat” and “zip_test.dat” files provided. Select the same 5% random sample (without replacement) of the train data used in exercise 1. Use this sample as your training data, and the complete test data for testing.
2. Define the response matrix (Y) and the predictor matrix (X). Center the predictor matrix.
3. Perform the Inter Batteries Analysis following the formulae given in the slides. Be aware that Y is not of full rank. Decide how many components you retain for prediction?.
4. Predict the responses in the test data, be aware of the appropriate centering. Compute the average R² in the test data.
5. Assign every test individual to the maximum response and compute the error rate. Compare the results with the obtained in exercise 1.