

Project: Thyroid disease

Task

Tell if a particular set of measurements (test results) comes from a person who is normal, or suffers from being hypothyroid or hyperthyroid (i.e. 3 output categories).

Data

There are 7200 observations representing patients. You are given 5000 of these, and 2200 are withheld for out-of-sample testing. There are 21 variables (there is no information on what these represent).

You are given the file *thyroidTrain.mat*, which contains the matrices *trainThyroidInput* (5000×21), *trainThyroidOutput* (5000×3), and *testThyroidInput* (2200×21). The first matrix, *trainThyroidInput*, contains the input patterns for the training data. The second matrix, *trainThyroidOutput*, contains the outputs coded in a “1-out-of-3” fashion. That is, the outputs are coded as (1,0,0), (0,1,0), or (0,0,1). The third matrix, *testThyroidInput*, contains the inputs for the test data. You are supposed to use the latter file to produce outputs that are handed in to me.

Steps and subgoals

1. Get acquainted with the data. Plot it and try to get a feel for the possible relationships between input and output.
2. Try to transform the variables and see if this changes the information content (this can be measured using e.g. a “Fisher Index”).
3. Construct a k -nearest neighbor (k -NN) classifier for the problem, using all the variables, and estimate the generalization error (use e.g. $k = 5$).
4. Prune the k -NN classifier by successively removing the variable that results in the least degradation of the generalization error, until the degradation is significant. Note the classification error (generalization).
7. Construct an artificial neural network (ANN) model using the remaining inputs. Optimize the number of hidden units (one hidden layer) with respect to the generalization error.
8. Try to prune the ANN model by successively removing the variable that results in the least degradation of the generalization performance, until the degradation is significant. Optimize the number of hidden units for the final model. Note the classification error.

9. Train a few networks with your optimal number of parameters. Combine these into a committee.
10. Hand in the test results for your best k -NN classifier and your best ANN committee together with your estimate of the generalization classification error.

Report and presentation of results

You will present the results from your project in two ways: (1) A written report where the main conclusions are presented together with figures and tables supporting your conclusions. (2) An oral presentation, of about 20 minutes, to your course colleagues.

The report should be about 10 pages, including figures and tables, and should contain the elementary report constituents:

- Introduction (brief presentation of problem, 1 page)
- Methodology (brief listing of methods, 1 page)
- Data (presentation of your data set with important observations, 1-2 pages)
- Results (4-5 pages)
- Discussion (your results and comparison to other researchers' results, 1 page)

The report writing should not take much more than one full day, since you are three persons sharing the work.

When you are finished with your report, and it has been accepted, then you should produce a postscript file with it, and pack it together with your dataset and other important parts of your project (like MATLAB M-files). The idea being that someone else could unpack it and repeat the main steps in your analysis without rewriting everything.