

SUPPLEMENTAL MATERIAL

CODUSA - Customize Optimal Donor Using Simulated Annealing In Heart Transplantation

D. Ansari, B. Andersson, M. Ohlsson, P. Höglund, R. Andersson, and J. Nilsson

Supplemental Methods

Survival analysis using neural networks

Cox regression models are the most commonly used methods for survival analysis in clinical research. Given its assumptions of proportionality of the hazard and the usual linear modelling of the covariates one would like to extend the analysis methods. Models using ANN have the ability to model the hazard with explicit time dependency and flexible nonlinear effects among the covariates. Biganzoli et al. showed that by treating the time interval as an input variable in a standard feed-forward ANN with a cross-entropy error function, it was possible to estimate smoothed discrete hazards as conditional probabilities of failure¹. The survival models used in this project follows the principles described in Biganzoli et al. with the extension of using ensembles of ANNs instead of a single one². For a general introduction to ANN see the Cross et al³. The ANNs were implemented as feed-forward multilayer perceptrons (MLP) with one hidden layer with the hyperbolic tangent as the activation function. The following error function was used during the training of the ANN models,

$$E = \sum_i \sum_l [d_{il} \log(h_l(x_i, a_l)) + (1 - d_{il}) \log(1 - h(x_i, h_l))]$$

where $h_l(x_i, a_l)$ is a smoothed estimate of the discrete hazard function for time interval l , which is modeled by the ANN output. The variables (x_i, a_l) represents the covariates for patient i and midpoint time for interval l , respectively. The event indicator d_{il} variable is one if uncensored patient i has the event in time interval l and zero otherwise. In order to have the possibility to regularize the ANN models a weight decay term was added to the above error function. This term, $E_r = \alpha \sum_j \omega_j^2$ introduces the parameter α , which is tuned during the model calibration procedure (see below). Finally to optimize the performance of the ANN model an ensemble approach were used, where several ANNs were combined into a single prediction model. The output of the ANN ensemble was computed as the mean of the output of the individual members in the ensemble. The ensemble was constructed by training the ANNs on different training sets, obtained from the random imputation technique when dealing with missing data. The ensemble size was 8 and no effort was used to optimize this

number. Given the discrete hazard function $h_l(x_i, a_l)$ and the definition $S(t_0) = 1$ the full survival curve can be constructed according to,

$$S(t) = \prod_{l:t_l < t} (1 - h_l)$$

where t_l is the end time for interval l .

Calibration and validation of the ANN models Calibration of each individual ANN was accomplished by minimizing the above error function using resilient back-propagation. To find the optimal regularization parameter and the optimal number of hidden nodes for the ANN 5-fold cross-validation was utilized. The number of hidden nodes was determined based on experiments starting with a single node and increasing the number of nodes until the highest accuracy was found for the validation sets. By a similar procedure the α -parameter was chosen to optimize the validation performance. When all parameters were set a new calibration using the full training dataset was performed. Throughout the model calibration the ensemble approach was used utilizing 8 different datasets, obtained from the missing data imputation technique (see below). The derivation cohort was used to calibrate and identify the optimal architecture for the ANN.

Risk variables identifications To identify important risk variables and to select the optimal set of risk variables used in the survival model, a ranking of risk variables was performed^{4,5}. A baseline C-index is created using all variables. The ranking list was then obtained by measuring the change of the C-index, as compared to the baseline, when a risk variable was excluded from the model. The highest ranked variable corresponds to the largest decrease of the C-index when it is excluded from the model. The lowest ranked variable will have the smallest effect on the C-index when excluded from the model and was subsequently be removed from the model. A new survival model was created and a new baseline C-index was computed, giving a new ranking list from which the lowest ranked variable again was removed. This backward elimination procedure was repeated until only one variable was left. The order in which the variables were removed constituted the final ranking list. Throughout this procedure full calibration of the model was performed, see Figure 1A for an illustration of the procedure including both model calibration and risk variable identification. The selection of the final set of variables was based in the obtained ranking list and was selected when no performance increase was found when adding the next variable from the ranking list.

Imputation of missing data

The ISHLT registry contains in average 29.3% missing data. We used the probability imputation technique multiple times to handle this^{6,7}. Each missing data was imputed 8 times with a random existing data point from another patient, which resulted in 8 study populations with a variation in variables that had missing data. Training and validation of the ANN models were performed with all of these 8 populations in an attempt to counterbalance random fluctuations and to utilize the ensemble approach.

Time dependent hazard ratio

The time dependent hazard ratios for the risk variables were determined in a similar way as described by Lippmann and co-workers⁸. By changing the risk variable in a patient from absent to present and calculating the hazard for the two conditions at each time interval, a time dependent hazard ratio for the specific risk variable of each patient could be determined. A hazard ratio for the specific variable was then obtained by computing the geometric mean of the hazard ratio from all patients. The 95% confidence intervals hazard ratio was calculated using the bootstrap technique (N= 10,000).

Software and computer resources

We performed the model building and simulation using high performance computer cluster with MatLab Distribution Computing Server 2010a, Neural Network Toolbox (MathWorks, Natick, Mass). We used the Lunarc Computational resources (www.lunarc.lu.se) and the Milleotto, an IBM bladecentre solution with 252 nodes containing two 64-bit, dual core (processors) Intel Xeon (3.0 GHz), corresponding to a total of 1,008 processors. An Apple Xserver® cluster (with 8 nodes, 64 cores) was also employed. Statistical analyses were performed using the Stata MP version 12.1 (2012) statistical package (StataCorp LP, College Station, Texas) and R version 2.15.1 (2012) and R version 2.15.1 (2012, The R Foundation for Statistical Computing).

Supplemental References

1. Biganzoli, E., Boracchi, P., Mariani, L., Marubini, E.. Feed forward neural networks for the analysis of censored survival data: a partial logistic regression approach. *Statistics in medicine* 1998;17(10):1169–86.
2. Hansen, L., Salamon, P.. Neural network ensembles. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 1990;12(10):993–1001.
3. Cross, S.S., Harrison, R.F., Kennedy, R.L.. Introduction to neural networks. *Lancet* 1995;346(8982):1075–9.
4. Nilsson, J., Ohlsson, M., Thulin, L., Nashef, S.A.M., Brandt, J.. Risk factor identification and mortality prediction in cardiac surgery using artificial neural networks. *The Journal of Thoracic and Cardiovascular Surgery* 2006;132:12–9.
5. Carlsson, A., Wingren, C., Kristensson, M., Rose, C., Fernö, M.r., Olsson, H.k., et al. Molecular serum portraits in patients with primary breast cancer predict the development of distant metastases. *Proceedings of the National Academy of Sciences of the United States of America* 2011;108:14252–7.
6. Schafer, J.L.. Multiple imputation: a primer. *Statistical methods in medical research* 1999;8(1):3–15.

7. Schemper, M., Smith, T.L.. Efficient evaluation of treatment effects in the presence of missing covariate values. *Statistics in medicine* 1990;9(7):777–84.
8. Lippmann, R.P., Shahian, D.M.. Coronary artery bypass risk prediction using neural networks. *The Annals of thoracic surgery* 1997;63(6):1635–43.