

Supplementary Online Content

Citterio D, Facciorusso A, Sposito C, Rota R, Bhoori S, Mazzaferro V. Hierarchic interaction of factors associated with liver decompensation after resection for hepatocellular carcinoma. *JAMA Surg*. Published online June 1, 2016. doi:10.1001/jamasurg.2016.1121.

eMethods. Recursive Partitioning Analysis

eTable. Two-Way Classification of Actual and Predicted LD Occurrence

eFigure. Ten-fold Cross-Validated Error Rate

This supplementary material has been provided by the authors to give readers additional information about their work.

eMethods. Recursive Partitioning Analysis

Recursive partitioning analysis is a statistical method that classifies patients in risk groups based on maximizing the value of the logistic regression (in case of classification trees) or of log-rank tests (in case of regression trees) for the clinical end-point of interest, by means of trees aimed at correctly classifying members of the population based on several independent variables¹ For selecting the splitting variable and its cutpoint, the model assesses all possible dichotomizations of all predictor variables to find the one dichotomization that produces the largest logistic regression test statistic. Following this, each split in the tree building process results in daughter nodes that are more “pure” than the parent node in the sense that groups of subjects with a majority for either response class are isolated. The method then repeats this assessment and in each of these two daughter nodes the variable that is most strongly associated with the response outcome (i.e. that produces the lowest p-value) is selected for the next split. In continuous splitting variables, the optimal cutpoint is also selected with respect to this criterion. This way, splitting continues recursively until some stop condition is reached. Finally, the terminal nodes are grouped according to their response outcome occurrence, and the results are presented as the final set of prognostic groups. The main feature of this model is that while in linear regression the information from different predictor variables is combined linearly, here the range of possible combinations includes all partitions that can be derived by means of recursive splitting, including multiple splits in the same variable². In ensemble methods such as the random forest (RF) a set of trees is built by individual trees originating from random samples of the original series. Usually, trees originated by ensemble methods are quite extended in absence of stopping or pruning rules. At the same time the random selection of each split of the tree may consider variables that would have otherwise been missed. While single classification trees are useful in identifying different risk classes – as the one represented in Figure 2 – ensembles of trees are not easy to interpret, because of lack in nesting capacity preventing the exact definition of classes. Ensemble methods are commonly used to validate and better define the importance of single parameters originated by means of recursive partitioning analysis, as they order different covariates in different contexts while identifying their potential effect on response. Ten-fold cross-validation refers to the process of dividing the original patient sample into 10 equal groups, then removing 1 group, used as validation sample, and reconstructing the model using the reduced sample set. The new model then is tested for predictive accuracy against the excluded fraction and the process is repeated 10 times (each time with a different excluded subset). Then, the average concordance index and error rate are calculated. This process is repeated 250 times to reduce the effect of random splits, and an overall c-index and error rate are calculated (Supplementary Figure 1).

eReferences

1. Breiman L, Friedman JH, Olshen RA, Stone CJ. *Classification and Regression Trees*. New York: Chapman and Hall; 1984.
2. Strobl C, Malley J, Tutz G. An introduction to recursive partitioning: rationale, application, and characteristics of classification and regression trees, bagging, and random forests. *Psychol Methods* 2009;14:323-348.

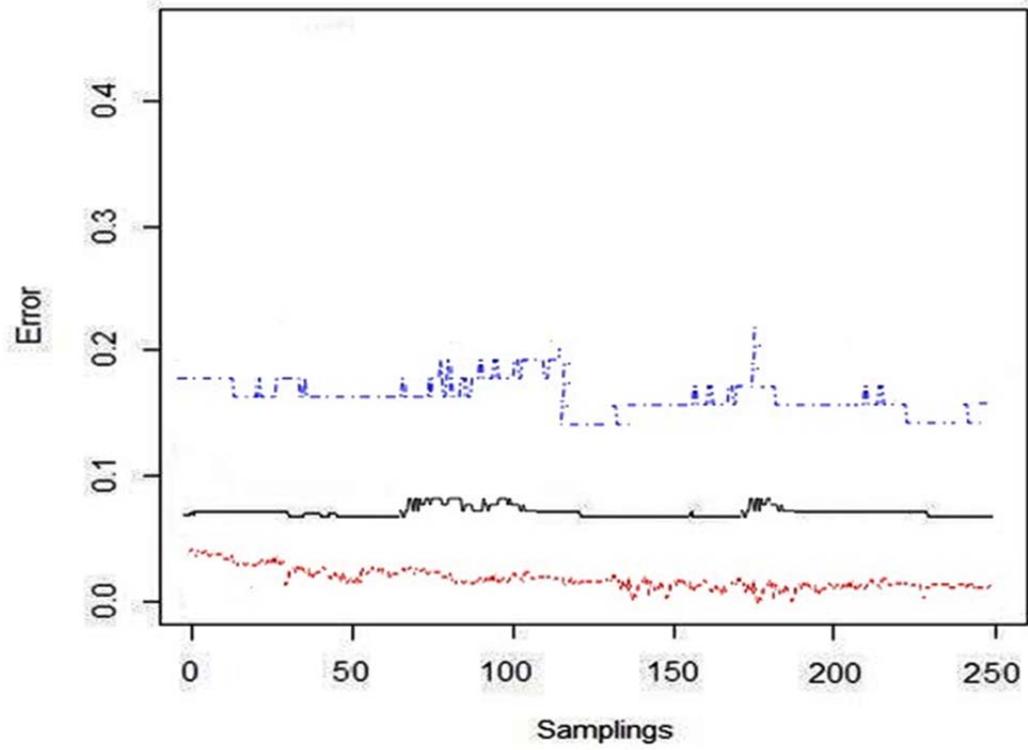
eTable. Two-Way Classification of Actual and Predicted LD Occurrence

Actual/Predicted	No LD	LD
No LD	421 (96.8%)	14 (3.2%)
LD	20 (18.5%)	88 (81.5%)
Error rate = 0.06		

LD indicates Liver Decompensation.

Error rate: the rate for which the response category with the highest fitted probability is not the observed category.

eFigure. Ten-fold Cross Validated Error Rate



The overall average error rate was 0.06. Blue line indicates error rate for LD prediction, red line indicates error rate for non-LD prediction, black line the average error rate.