



Exploring Decision Trees with R

J. Alfredo Freites

OCRUG, February 2018



Sources

rpart *vignette* <http://127.0.0.1:30960/help/library/rpart/doc/longintro.pdf>

rpart.plot *vignette* <http://127.0.0.1:30960/help/library/rpart.plot/doc/prp.pdf>

**Max Kuhn and Kjell Johnson *Applied
Predictive Modeling*, Springer, 2013. ISBN
978-1-4614-6848-6**

**Tom Mitchell *Machine Learning*, McGraw-
Hill, 1997. ISBN 0070428077**



Decision Tree modeling

Partition the data into smaller groups that are more homogeneous with respect to the response

real-valued response: Regression Trees

categorical (discrete): Classification Trees

Known knowns (the bad)

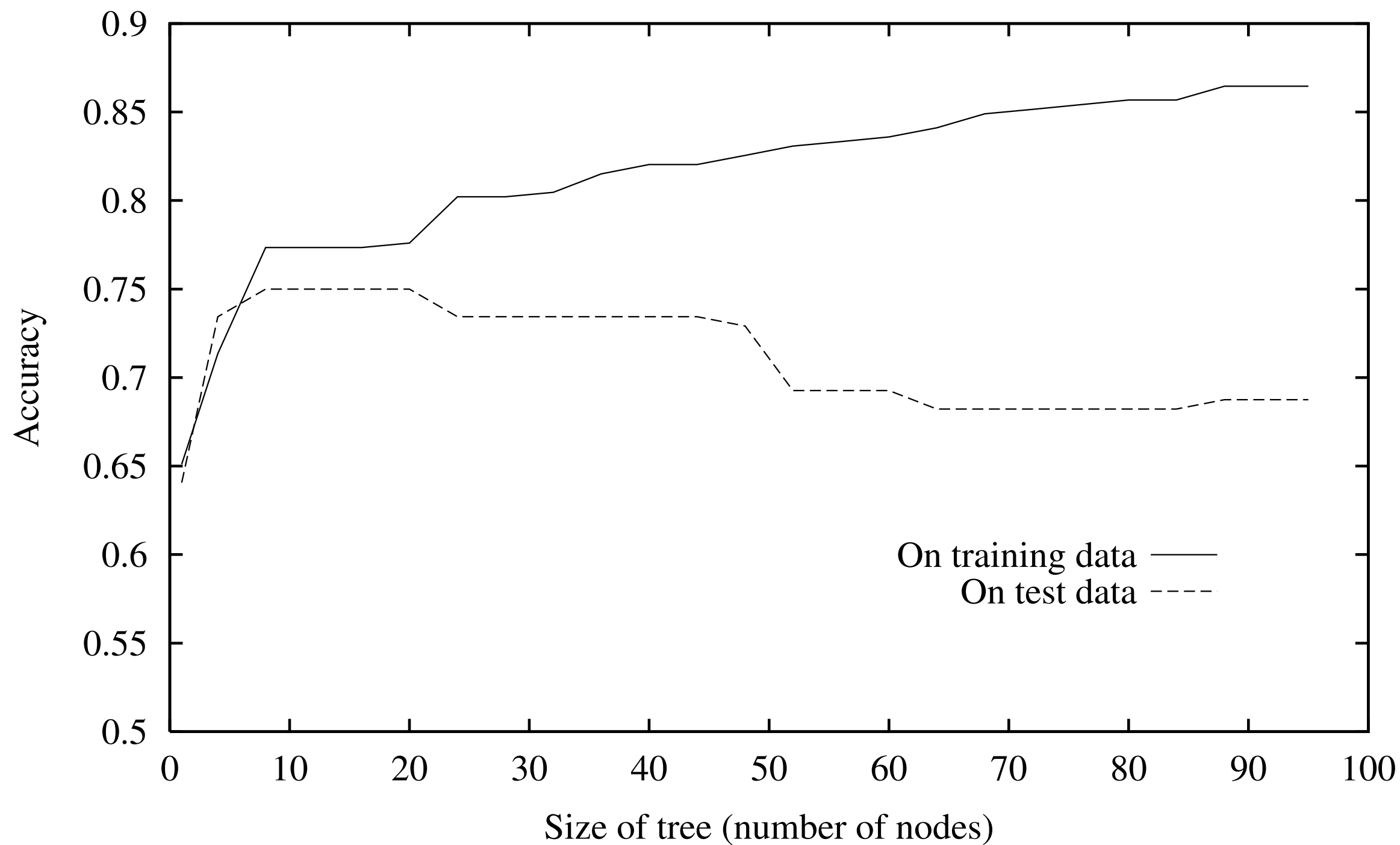
Overfitting

Model instability

If-then statements produce rectangular partitions, which may negatively affect predictive performance



Overfitting





Decision Tree modeling

Known knowns (the good)

White box. Highly interpretable. I know how it was done.

Can handle all kinds of variables.

Can handle missing data.

Base for successful ensemble methods.



rpart (recursive partitioning)

CART (Classification and Regression Trees)

Breiman, Friedman, Olshen and Stone book.

can do

Categorical

Linear regression

Poisson

Exponential



Two stages

First:

Find the single variable which BEST splits the data into two groups

***Apply separately to each sub-group
RECURSIVELY until minimum size or no improvement***

Second:

Use cross-validation to prune the resulting tree



Key points

Splitting criteria: Gini, entropy.

Pruning method: Cross validation

Missing data: if there's at least one independent variable, it will be used!



Maximize this thing

I: Impurity measure

p: probability

A: the parent node (a rule and a set of data)

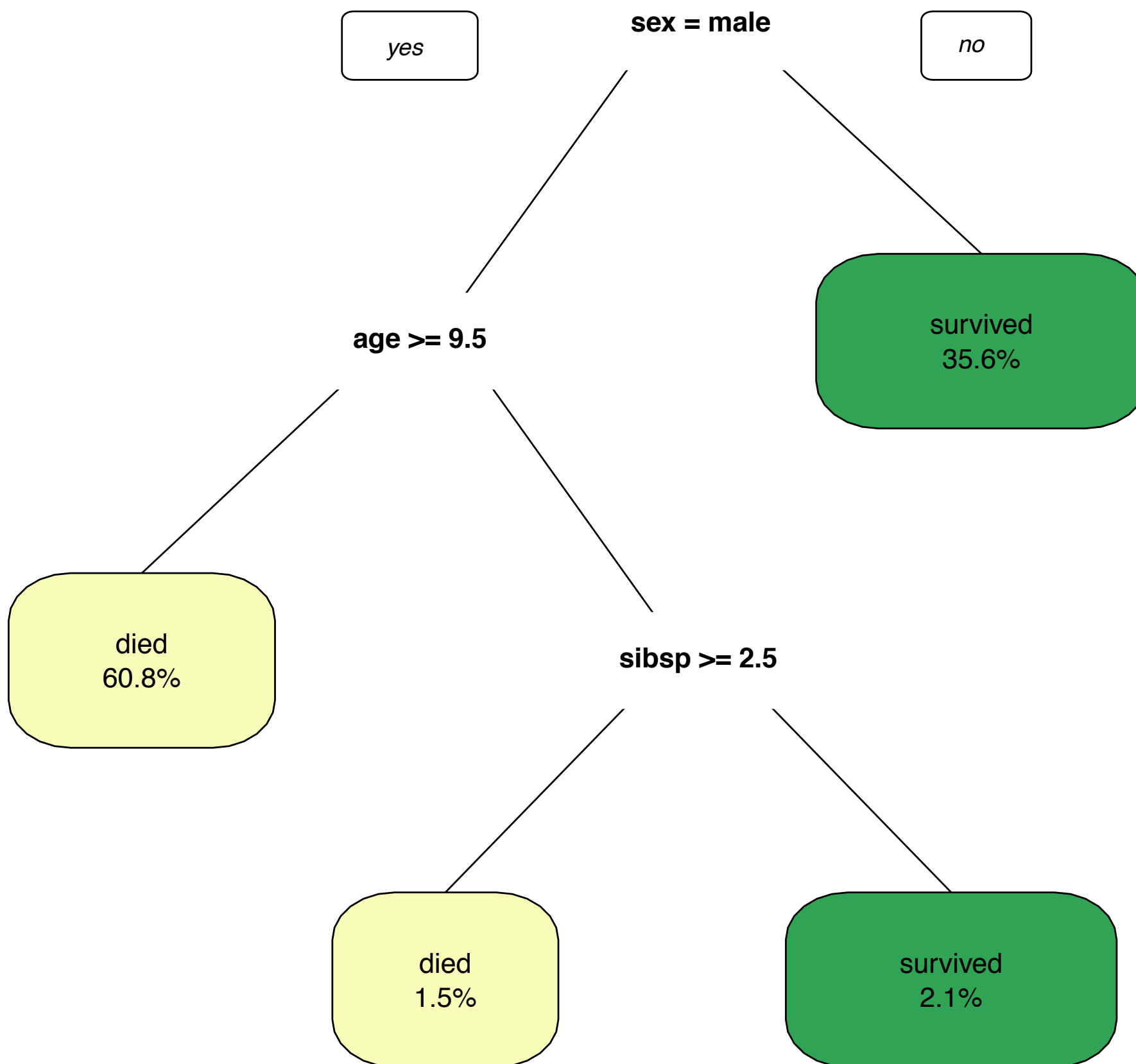
A_L, A_R: the left and right children

$$\Delta I = p(A)I(A) - p(A_L)I(A_L) - p(A_R)I(A_R)$$



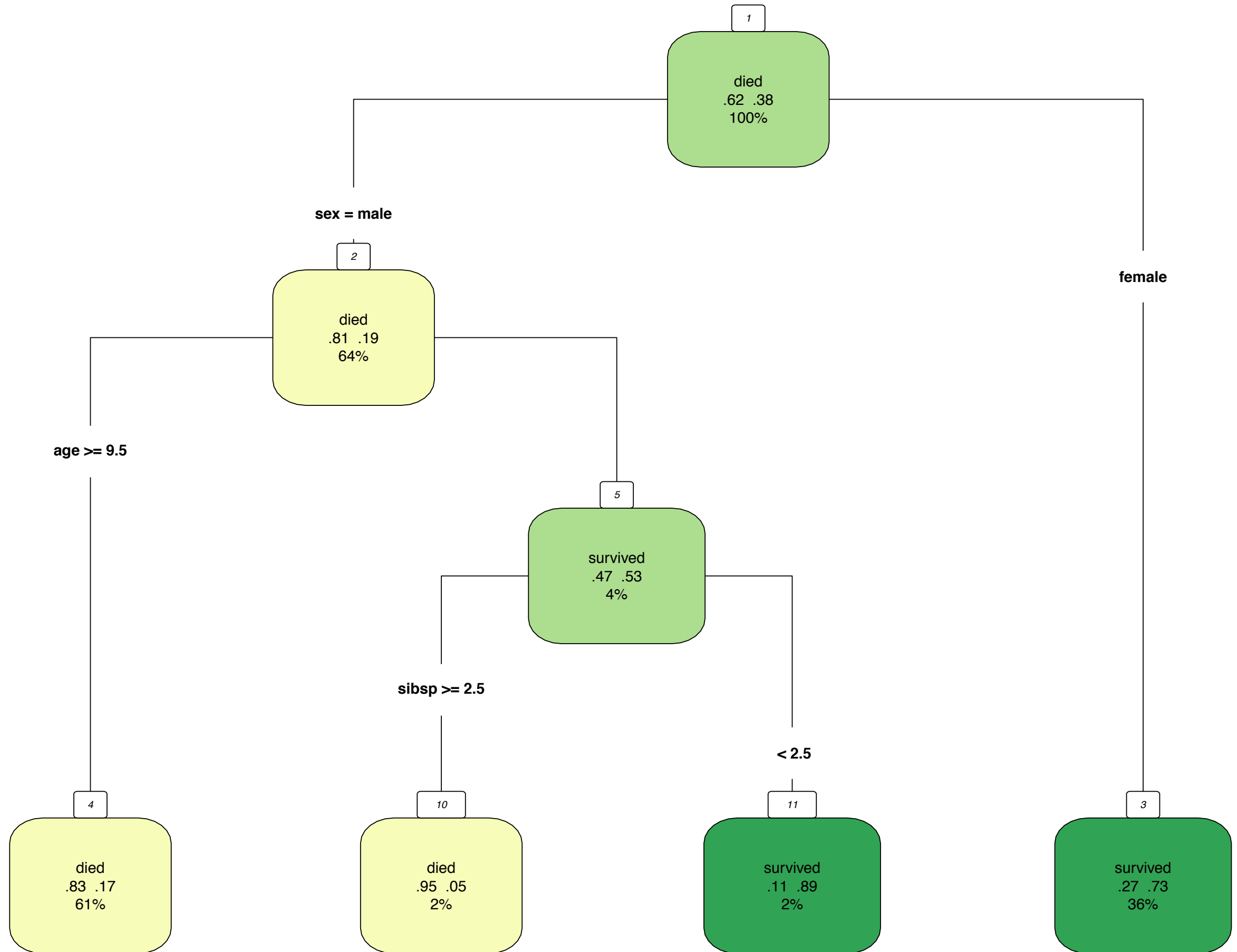


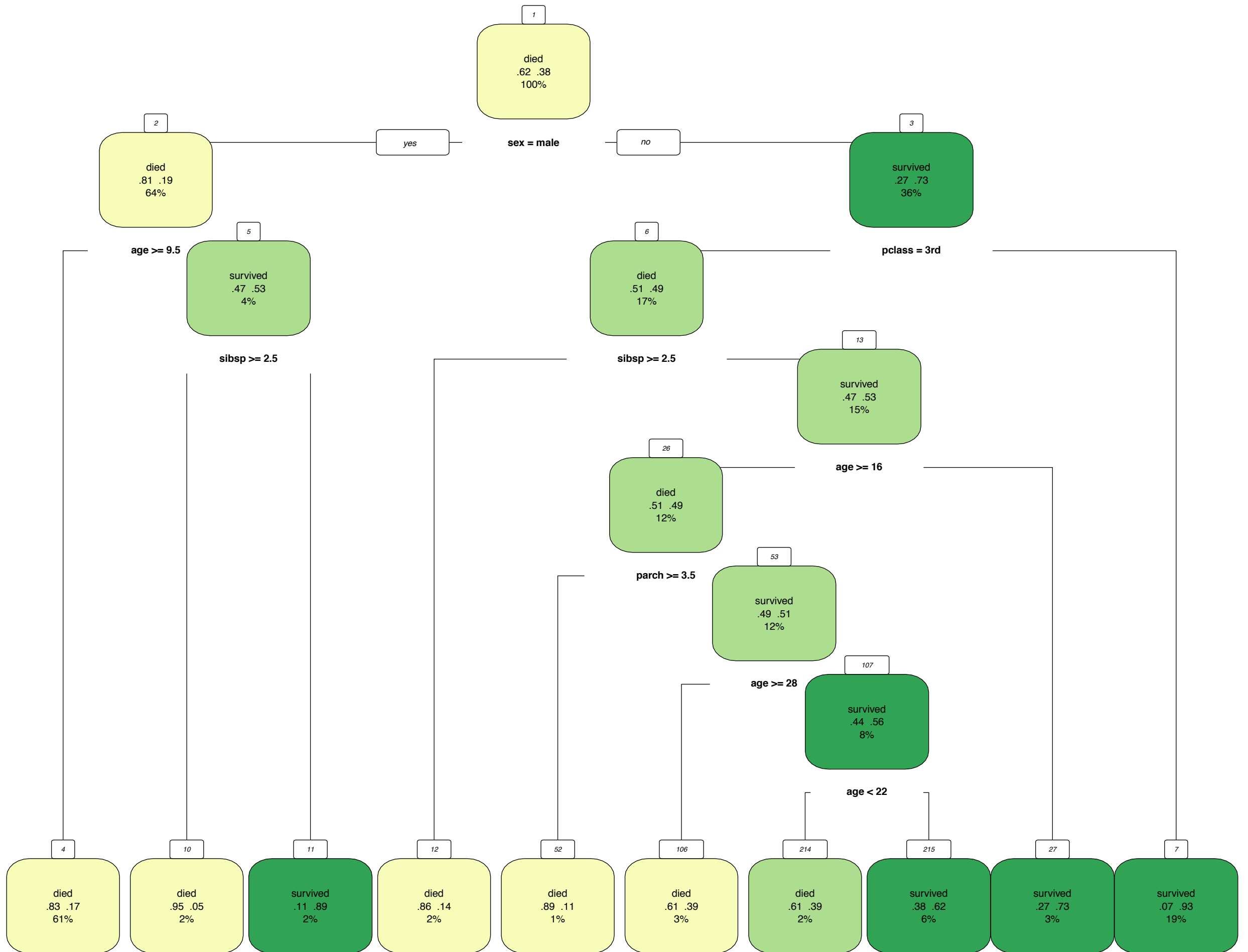
Who survived?





Who Survived?







Cost of splitting

cp: complexity parameter

R is the risk

T is the number of terminal nodes

T₁ is the one-node tree

$$R_{cp}(T) \equiv R(T) + cp * |T| * R(T_1)$$