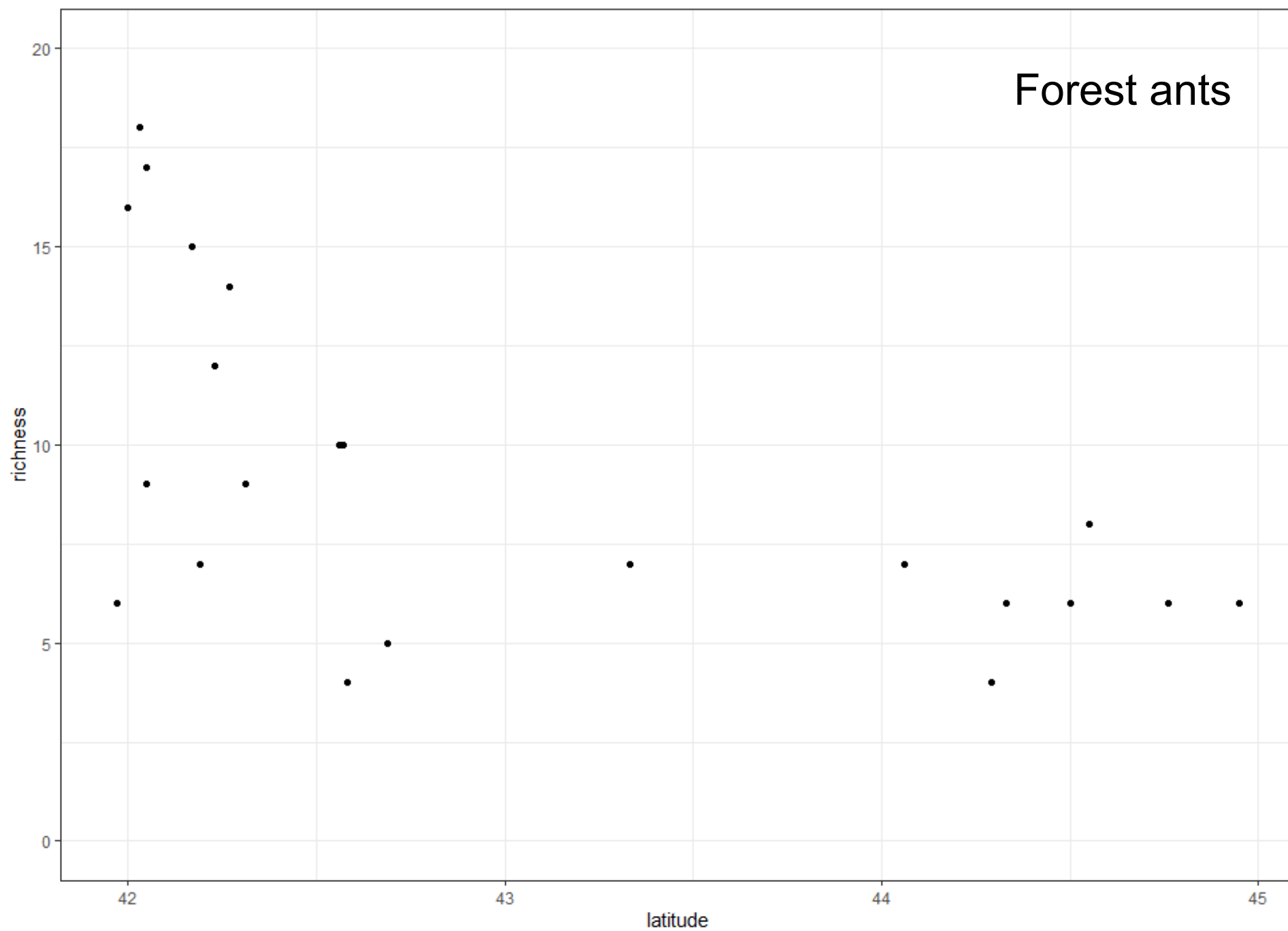
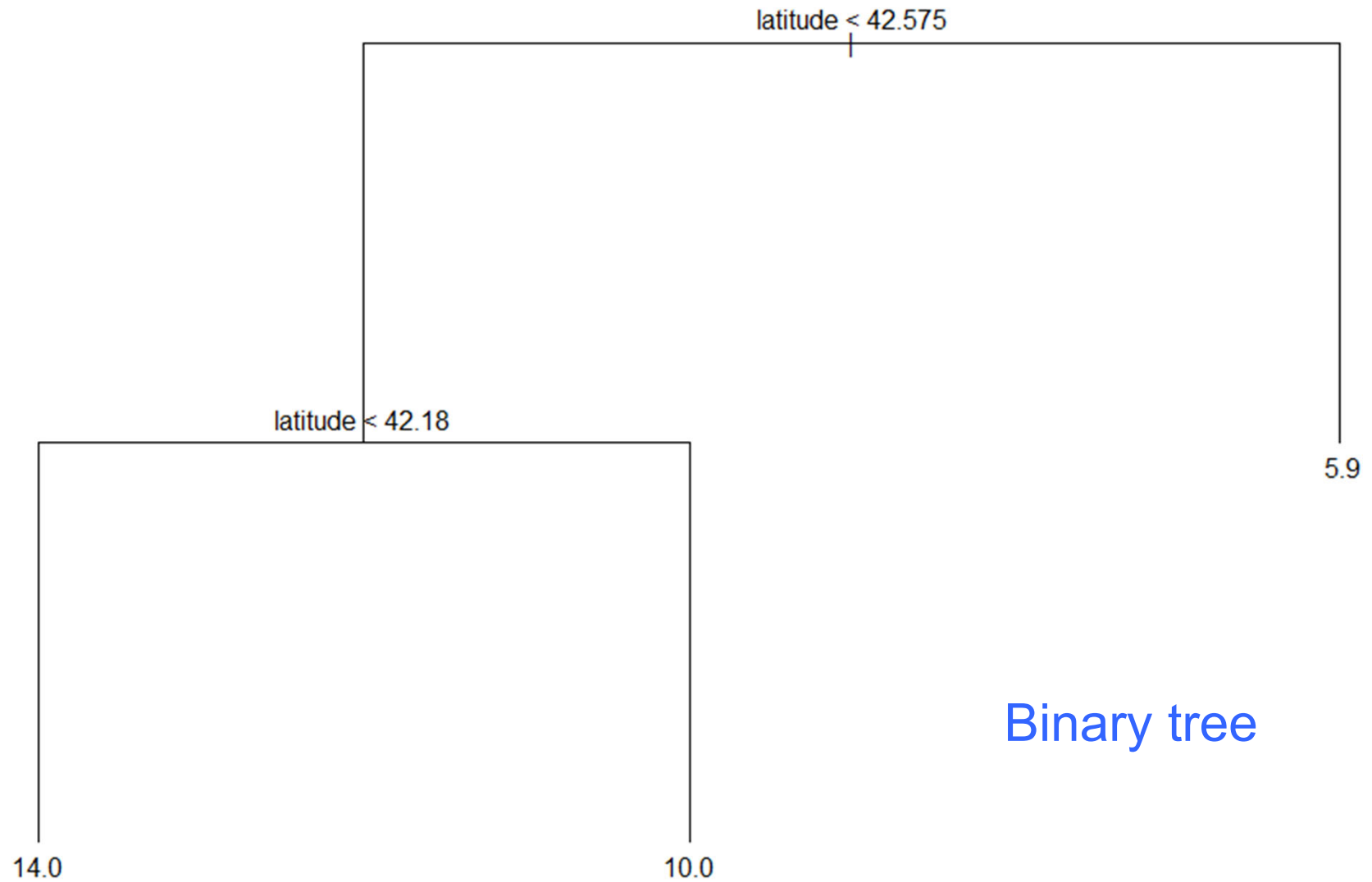


Today

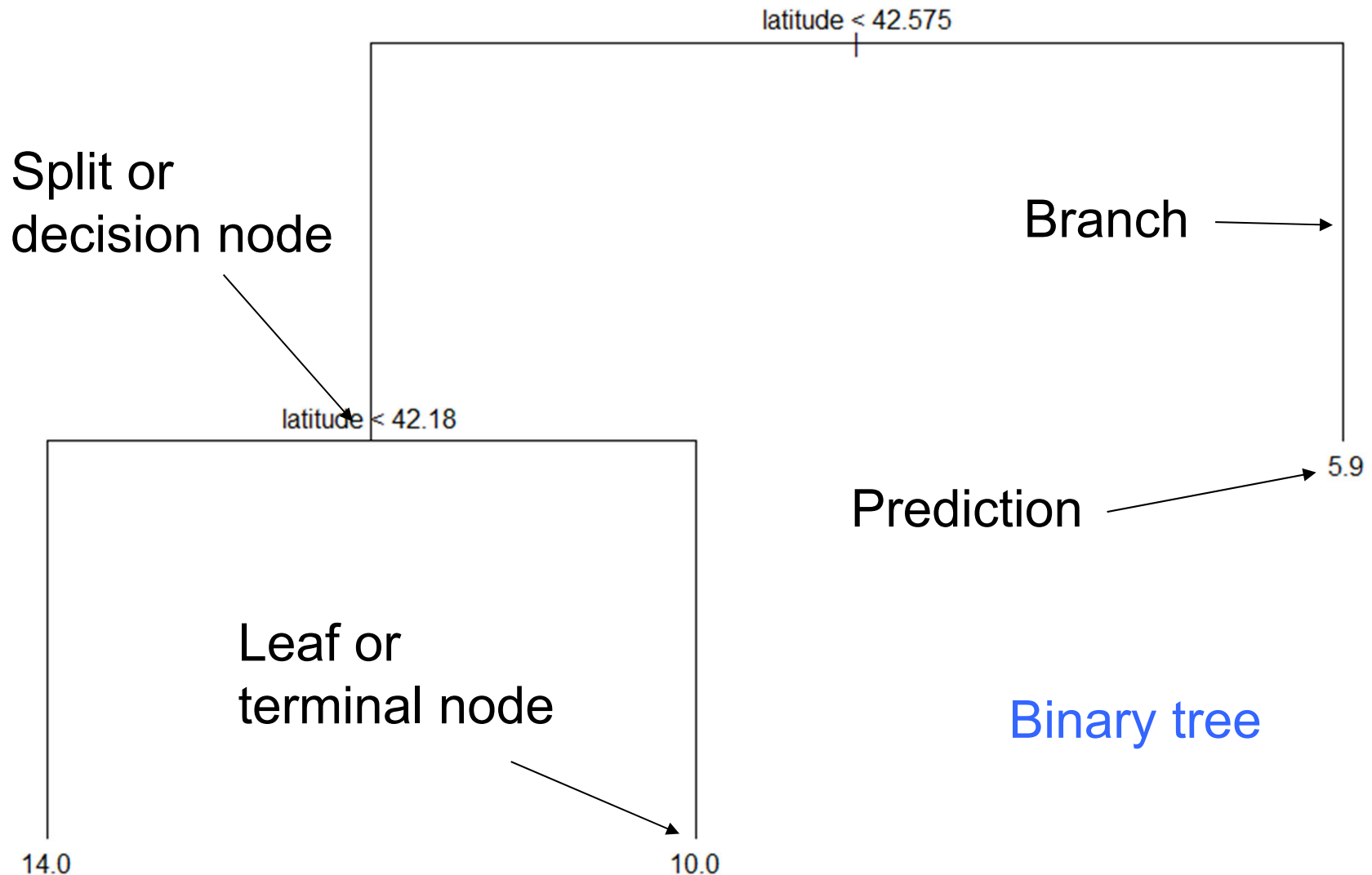
- Decision tree models
- + training and inference



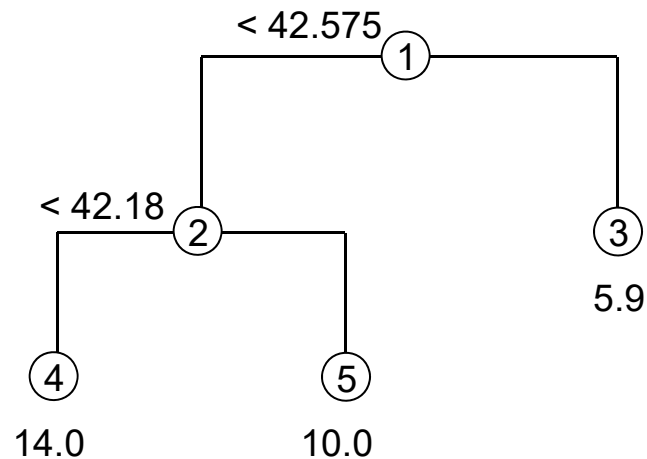
A regression tree model



A regression tree model



Model algorithm

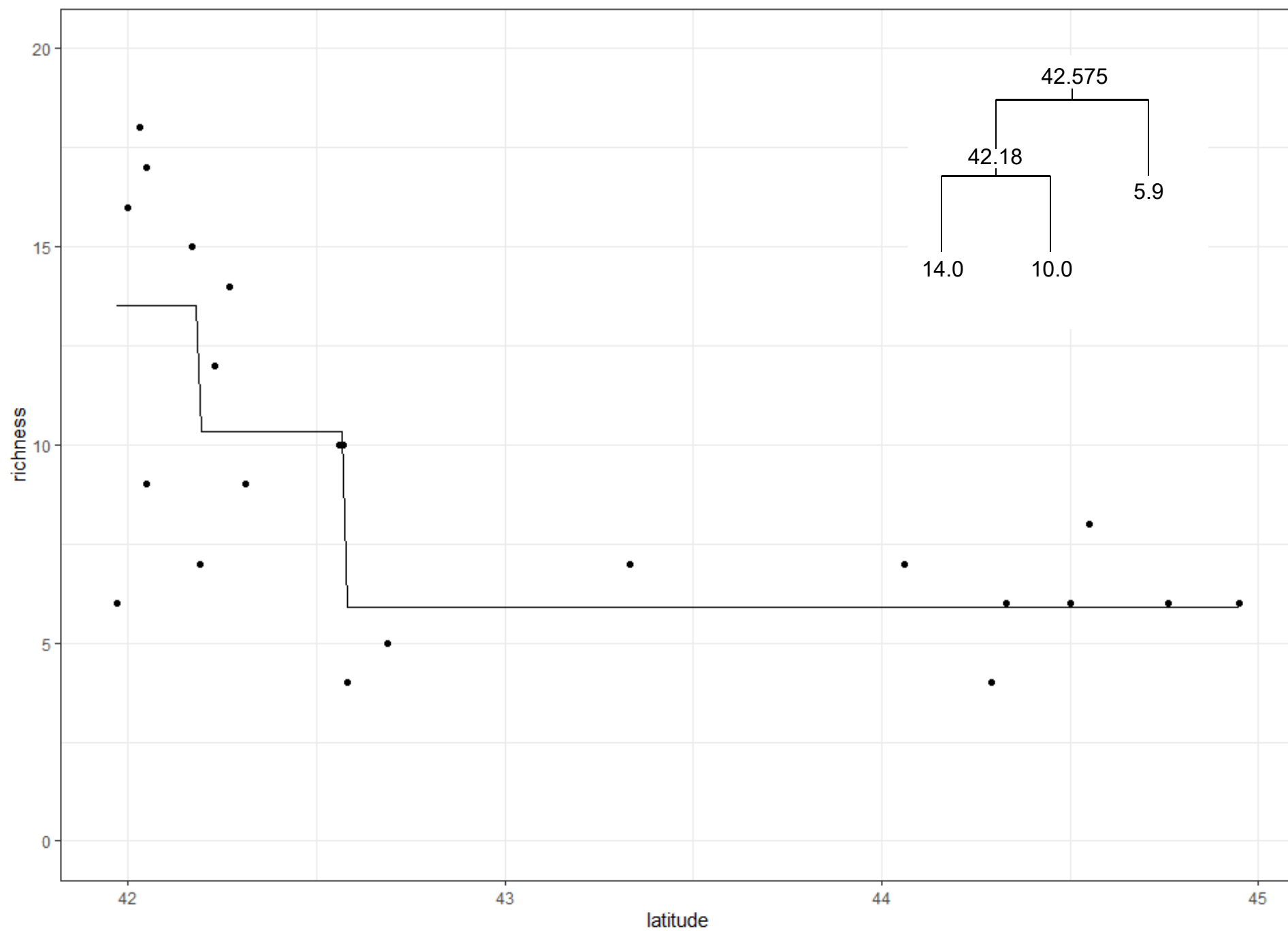


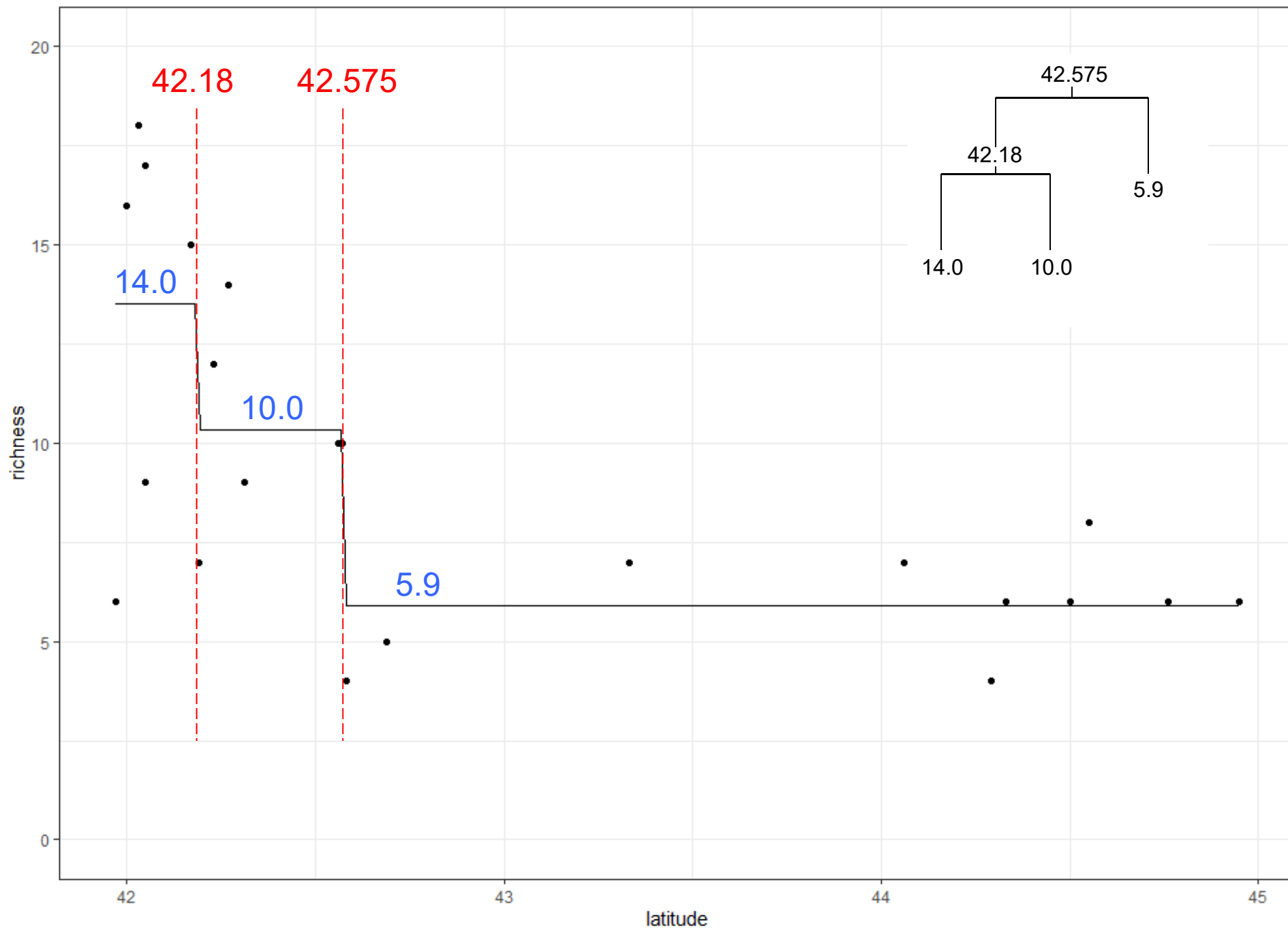
Tree structure

node	type	split	y
1	split	42.575	
2	split	42.180	
3	leaf		5.9
4	leaf		10.0
5	leaf		14.0

Algorithm

```
start at the root node
while node type is split
  if x < split value at the node
    take left branch to next node
  else
    take right branch to next node
return predicted y at node
```





Training algorithm

Binary recursive partitioning

```
define build_tree(y, x)
  if stop = TRUE
    calculate prediction (mean of y)
  else
    find x_split  #best x to split the data
    build_tree( (y, x)[x < x_split] )  #L branch
    build_tree( (y, x)[x >= x_split] )  #R branch
```

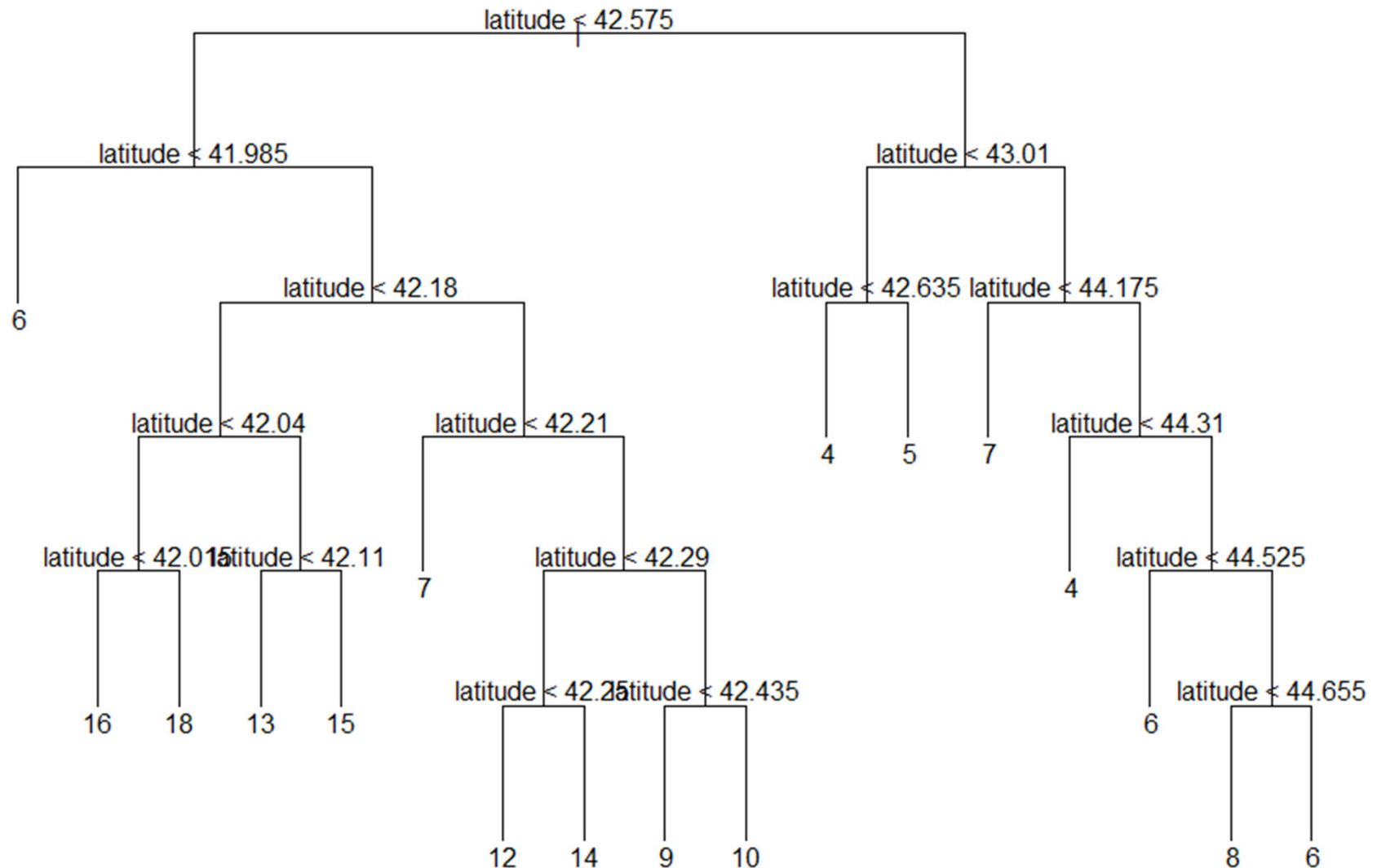
Stopping rules e.g.

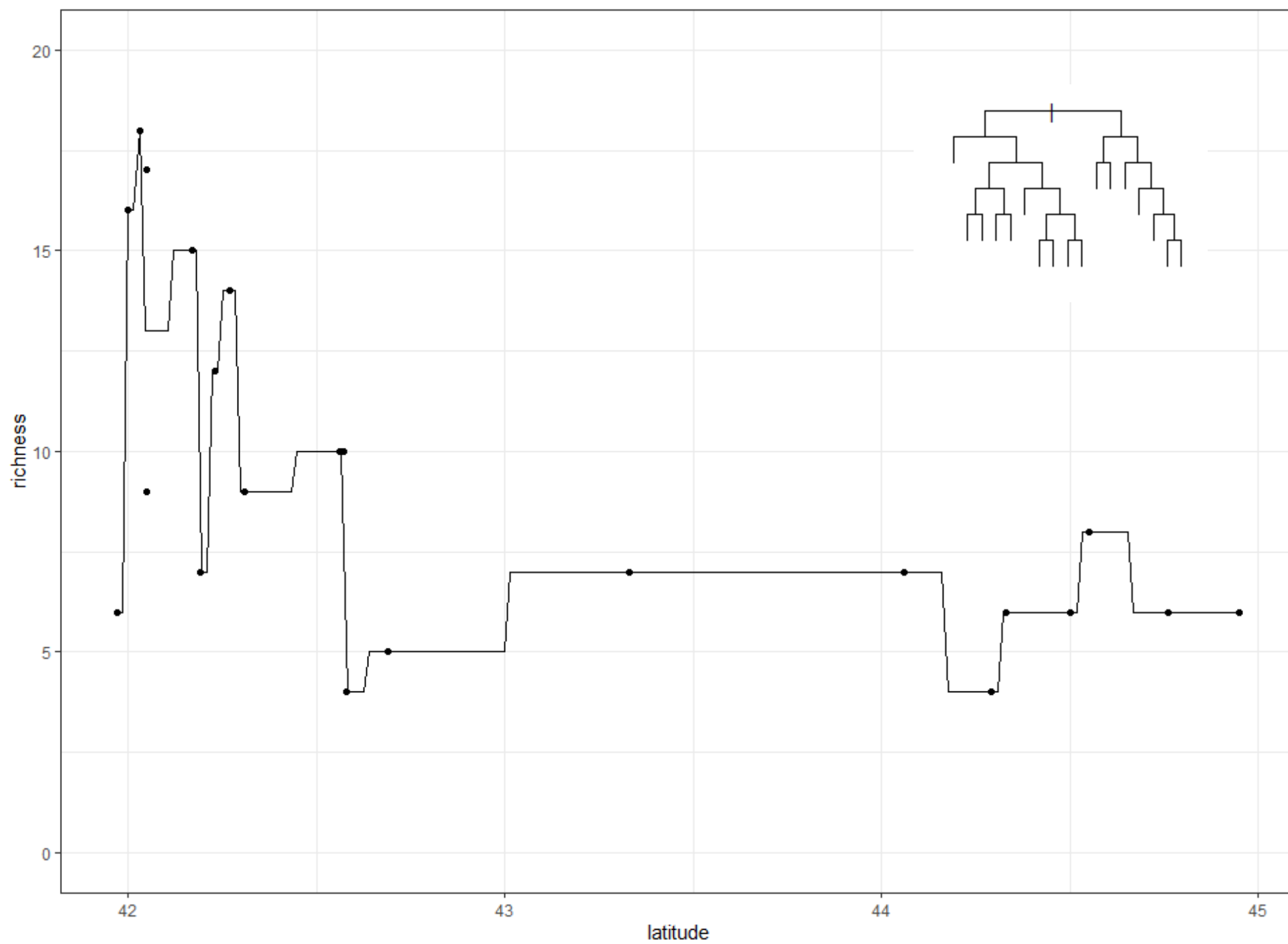
- data per node
- tree depth
- node variance
- error improvement

Find splits that minimize training error

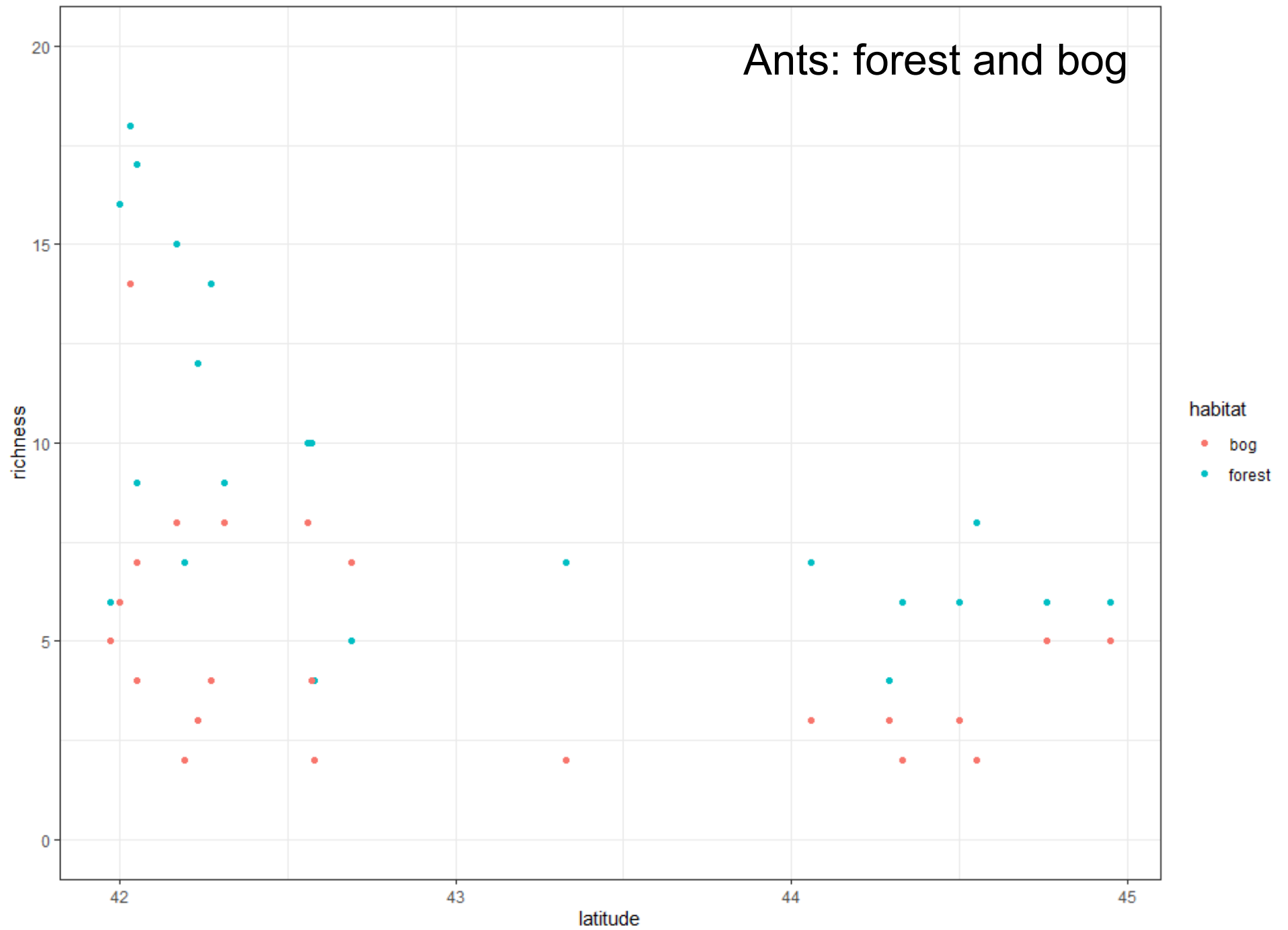
- regression: SSQ
- classification: Gini index or entropy

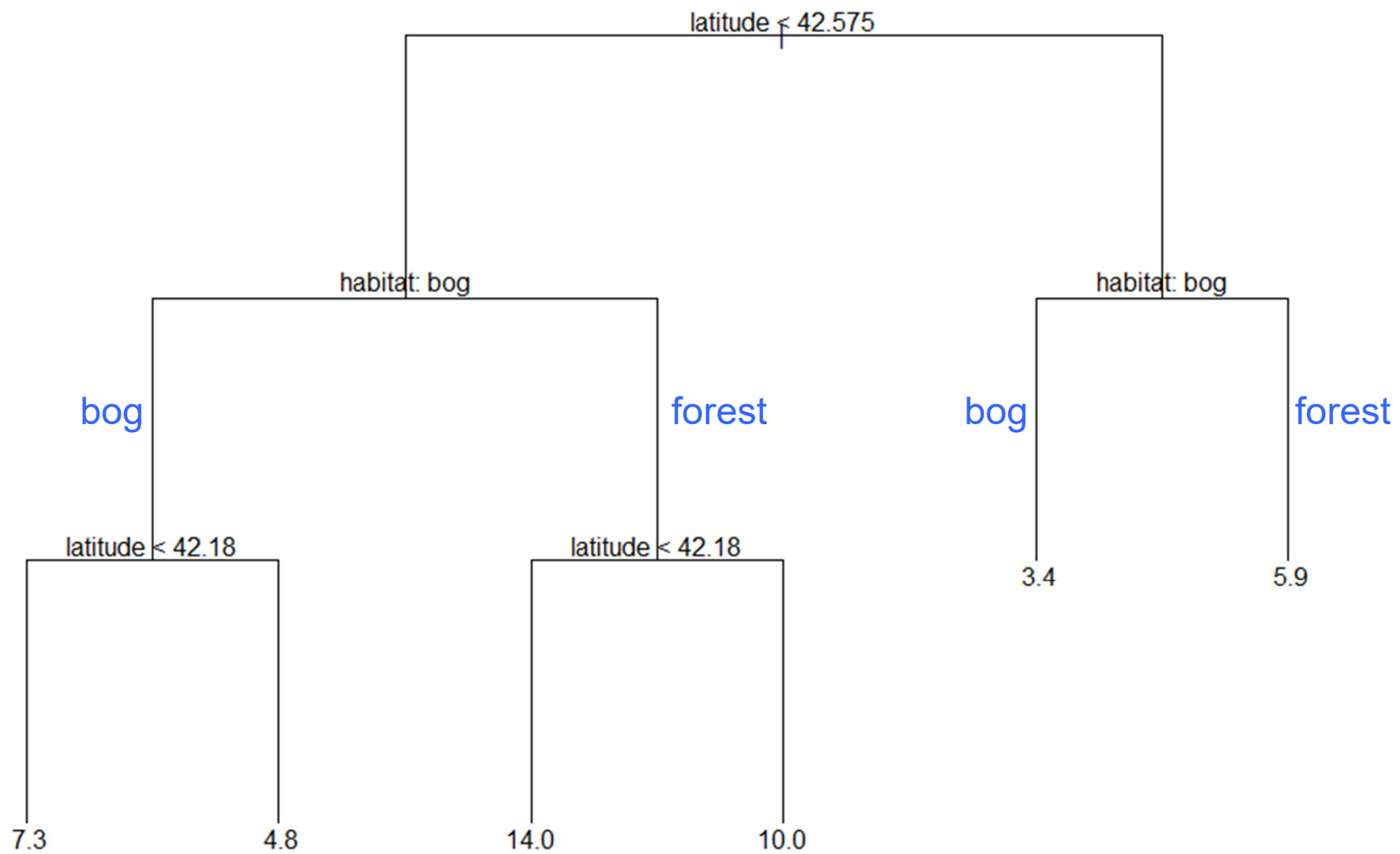
Same data, deeper tree



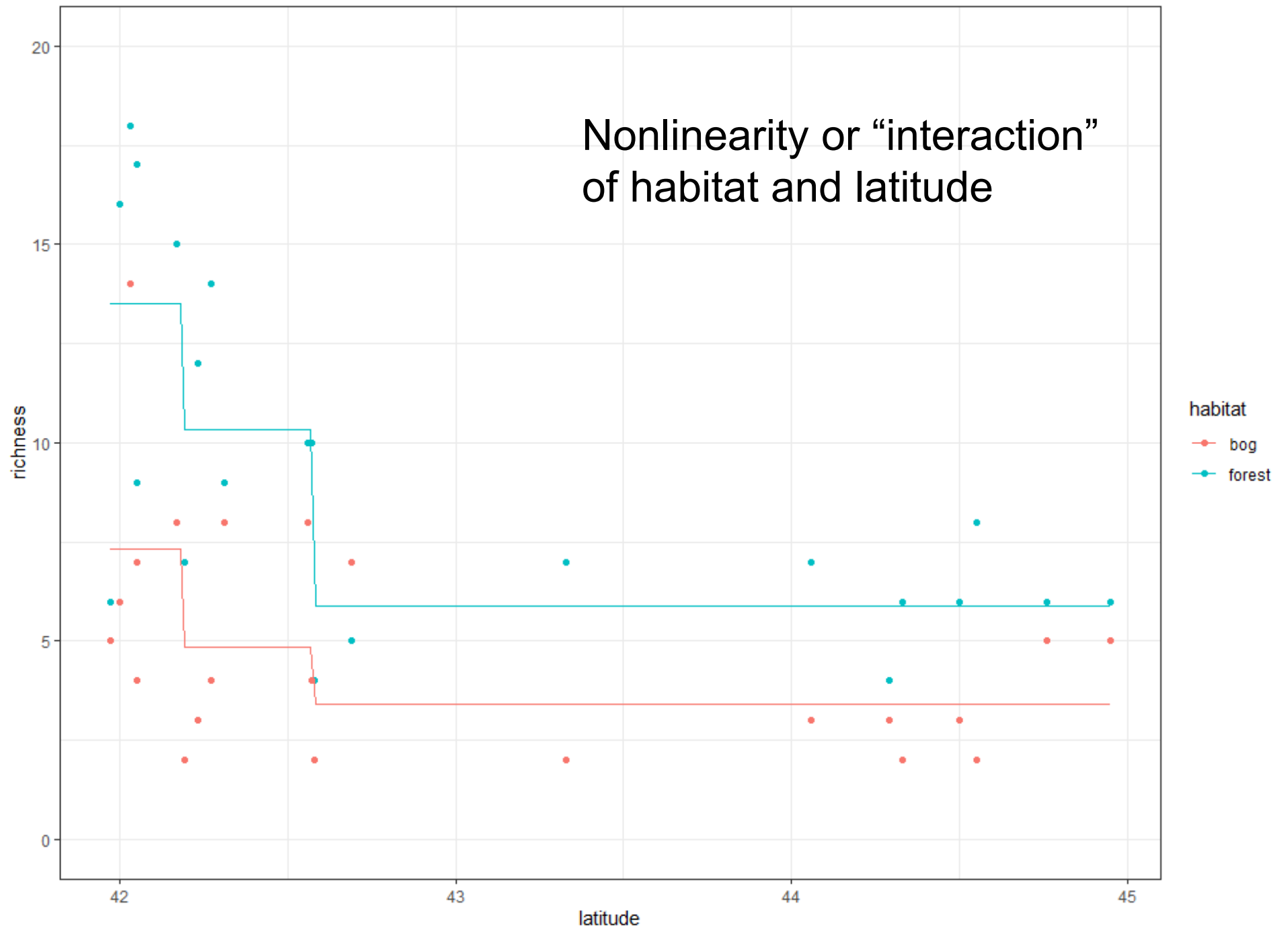


Ants: forest and bog



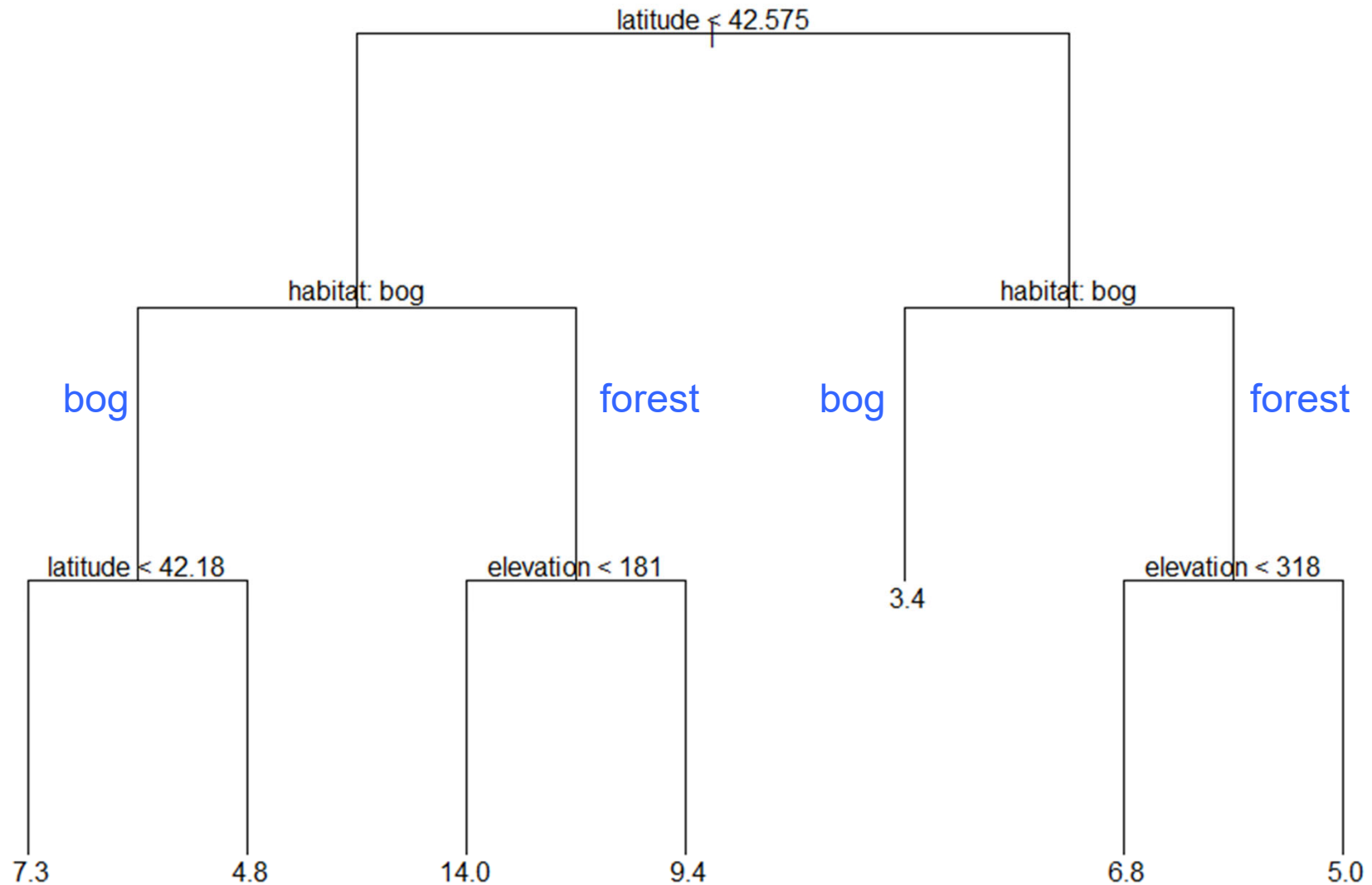


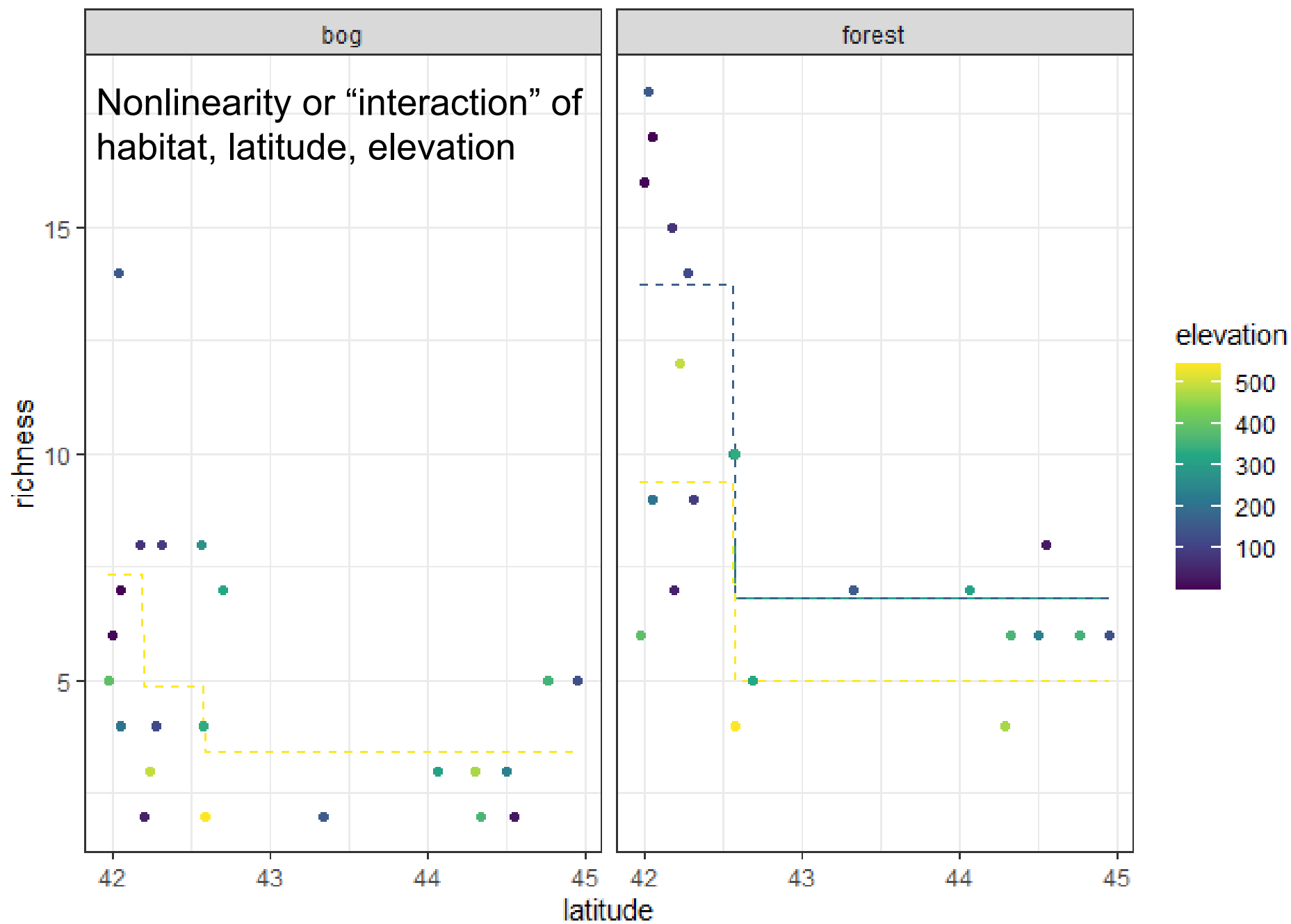
Nonlinearity or “interaction” of habitat and latitude



```
> head(ants)
  habitat latitude elevation richness
1 forest    41.97      389         6
2 forest    42.00         8        16
3 forest    42.03     152        18
4 forest    42.05         1        17
5 forest    42.05     210         9
6 forest    42.17         78        15
```

All 3 predictors





Inference

- k-fold CV
- Can also use for tree complexity
 - training: complexity penalty
 - e.g. $\text{loss} = \text{SSQ} + \alpha T$
 - where α is a tuning parameter, T is number of leaves
 - “pruning”

