

# Using CV for assessment (I)

→ CROSS VALIDATION

How the learned artifact will behave on unseen data?

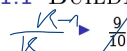
More precisely:

How an artifact learned with **this learning technique** will behave on unseen data?

## Using CV for assessment (II)

“This learning technique” = `BUILDDECISIONTREE()` with  $k_{\min} = 10$

1. repeat  $k$  times

1.1 `BUILDDECISIONTREE()` with  $k_{\min} = 10$  on all but one slice  


1.2 compute classification error on left out slice

2. average computed classification errors

  
10 invocations of `BUILDDECISIONTREE()`

## Using CV for assessment (III)

“This learning technique” = `BUILDDECISIONTREE()` with  $k_{\min}$  chosen automatically with a 10-fold CV

For assessing this technique, we do two nested CVs:

1. repeat  $k$  times
  - 1.1 choose  $k_{\min}$  among  $m$  values with 10-CV (repeat `BUILDDECISIONTREE()`  $10m$  times) on all but one slice
    - ▶  $\frac{k-1}{k} \frac{9}{10} n$  observations in each  $\mathbf{X}$  passed to `BUILDDECISIONTREE()`!
  - 1.2 compute classification error on left out slice
    - ▶ usually, a new tree is built on  $\frac{k}{k-1} n$  observations
2. average computed classification errors

$(10 + 1)k$  invocations of `BUILDDECISIONTREE()`

## Using CV for assessment: “cheating”

“This learning technique” = `BUILDDECISIONTREE()` with  $k_{\min}$  chosen automatically with a 10-fold CV

Using just one CV is cheating (cherry picking)!

- ▶  $k_{\min}$  is chosen exactly to minimize error on the full dataset
- ▶ conceptually, this way of “fitting”  $k_{\min}$  is similar to the way we build the tree

## Subsection 1

### Regression trees

# Regression with trees

Trees can be used for regression, instead of classification.

decision tree vs. regression tree

## Tree building: decision $\rightarrow$ regression

**function** BUILDDECISIONTREE( $\mathbf{X}, \mathbf{y}$ )

**if** SHOULDSTOP( $\mathbf{y}$ ) **then**

$\hat{y} \leftarrow$  most common class in  $\mathbf{y}$

**return** new terminal node with  $\hat{y}$

**else**

$(i, t) \leftarrow$  BESTBRANCH( $\mathbf{X}, \mathbf{y}$ )

$n \leftarrow$  new branch node with  $(i, t)$

append child BUILDDECISIONTREE( $\mathbf{X}|_{\mathbf{x}_i < t}, \mathbf{y}|_{\mathbf{x}_i < t}$ ) to  $n$

append child BUILDDECISIONTREE( $\mathbf{X}|_{\mathbf{x}_i \geq t}, \mathbf{y}|_{\mathbf{x}_i \geq t}$ ) to  $n$

**return**  $n$

**end if**

**end function**

**Q:** what should we change?

$\in \mathbb{R}^m$

# Tree building: decision $\rightarrow$ regression

```
function BUILDDECISIONTREE( $\mathbf{X}, \mathbf{y}$ )  
  if SHOULDSTOP( $\mathbf{y}$ ) then  
     $\hat{y} \leftarrow \bar{y}$  ▷ mean  $\mathbf{y}$   
    return new terminal node with  $\hat{y}$   
  else  
     $(i, t) \leftarrow \text{BESTBRANCH}(\mathbf{X}, \mathbf{y})$   
     $n \leftarrow$  new branch node with  $(i, t)$   
    append child BUILDDECISIONTREE( $\mathbf{X}|_{\mathbf{x}_i < t}, \mathbf{y}|_{\mathbf{x}_i < t}$ ) to  $n$   
    append child BUILDDECISIONTREE( $\mathbf{X}|_{\mathbf{x}_i \geq t}, \mathbf{y}|_{\mathbf{x}_i \geq t}$ ) to  $n$   
    return  $n$   
  end if  
end function
```

**Q:** what should we change?



## Best branch

```
function BESTBRANCH(X, y)  
     $(i^*, t^*) \leftarrow \arg \min_{i,t} E(\mathbf{y} | \mathbf{x}_i \geq t) + E(\mathbf{y} | \mathbf{x}_i < t)$   
    return  $(i^*, t^*)$   
end function
```

**Q:** what should we change?

# Best branch

**function** BESTBRANCH(**X**, **y**)

$(i^*, t^*) \leftarrow \arg \min_{i,t} \sum_{y_i \in \mathbf{y} | x_i \geq t} (y_i - \bar{y})^2 + \sum_{y_i \in \mathbf{y} | x_i < t} (y_i - \bar{y})^2$

**return**  $(i^*, t^*)$

**end function**

**Q:** what should we change?

Minimize sum of residual sum of squares (RSS) (the two  $\bar{y}$  are different)

DIFFERENT

elements of  $\mathbf{y}$

label of decision node

## Stopping criterion

```
function SHOULDSTOP(y)  
  if y contains only one class then  
    return true  
  else if  $|\mathbf{y}| < k_{\min}$  then  
    return true  
  else  
    return false  
  end if  
end function
```

**Q:** what should we change?

## Stopping criterion

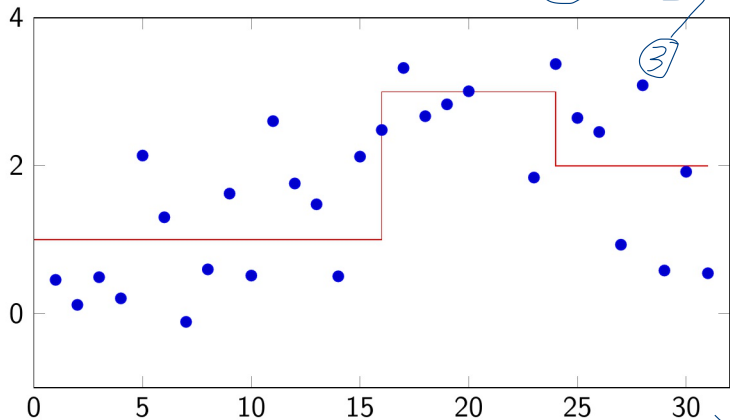
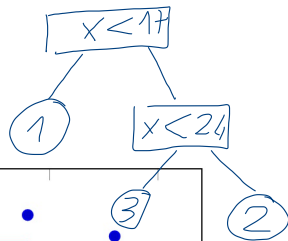
```
function SHOULDSTOP(y)
  if RSS is 0 then
    return true → RARE
  else if  $|y| < k_{\min}$  then
    return true
  else
    return false
  end if
end function
```

Q: what should we change?

# Interpretation

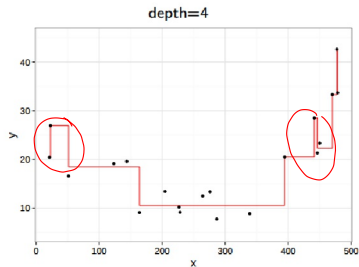
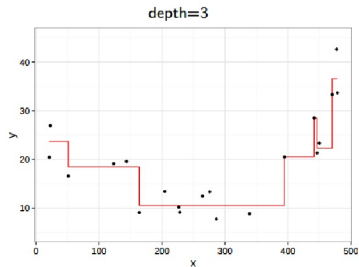
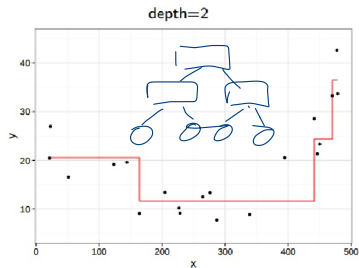
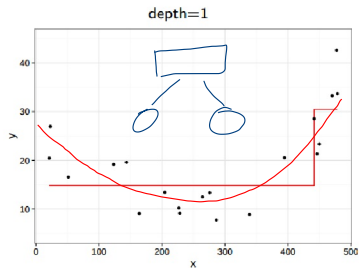
$$p = 1$$
$$n = 29$$

y



x

# Regression and overfitting



# Trees in summary

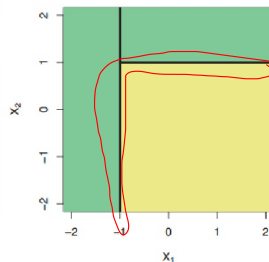
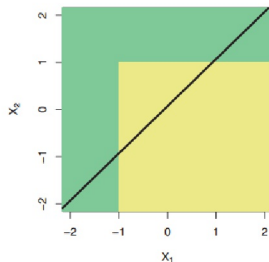
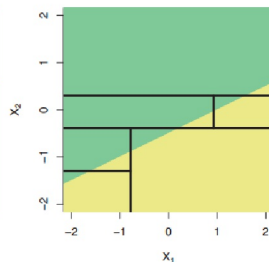
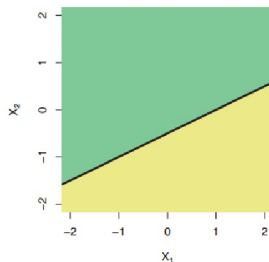
## Pros:

- ▲ easily interpretable/explicable
- ▲ learning and regression/classification easily understandable
- ▲ can handle both numeric and categorical values

## Cons:

- ▼ not so accurate (**Q**: always?)

# Tree accuracy?



DECISION  
BOUNDARY



## Lab: tree on iris (2 h)

- ▶ for each of the 5 variables in iris, predict it with the other 4
- ▶ which is the hardest to be predicted? why?

Packages: tree

Functions: tree, prune.tree, `predict.tree(t, type="class")`, `table`

CONFUSION  
MATRIX

→ LEARNING

$X, y$        $\text{Species} \sim$   
formula, d