

① Try all split rules. Calc  $G_{int}$  for each node. Select split rule with lowest weighted  $G_{int}$  avg:

$$\frac{h_L G_L + h_R G_R}{h_L + h_R}$$

② Repeat step 1 recursively until nodes has  $N_0$ . Default  $N_0 = 1$ .  
Assign  $\hat{y} = \text{Mode}\{\text{y in node}\}$

How to pick  $N_0$ ? *arbitrary!*

RF? Yes!

Don't

Don't

For  $N_0 = 1$  because splitting is good is *begging!*

Is *mischg.* even *re* *auf* *gerade*? No...

Now topic: asymmetric cost models.

There are two errors; wrongs s.t.  $\hat{y} \neq y$

① If  $\hat{y} = 0$  when  $y = 1 \Rightarrow$  False negative (FN)

② If  $\hat{y} = 1$  when  $y = 0 \Rightarrow$  False positive (FP)

all *gerade*  
can be summed up in a  $2 \times 2$  table called a confusion table!

		$\hat{y}$		
		0	1	
$y$	0	TN	FN	#N
	1	FN	TP	#P
		#FN	#TP	n

from regression  
true pos  
num neg.  
num pos.  
sample size

num  
pred  
key

num  
pred.  
pos.

$$err := \frac{FP + FN}{n}$$

$$Acc = \frac{TP + TN}{n} = 1 - \frac{FP + FN}{n} = 1 - err$$

other  
many metrics

$$precision = \frac{TP}{\#PP}$$

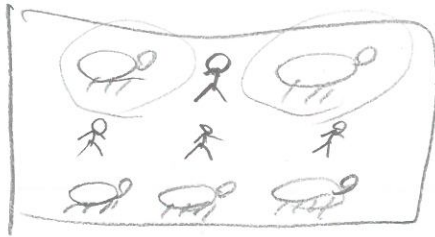
what prop of persons did you locate?

$$recall = \frac{TP}{\#P}$$

what prop of true  
pos. were actually pos?

Task: identify bugs

12



$$\text{Precision} = \frac{\#TP}{\#PP} = \frac{2}{2} = 100\% \quad \text{good!}$$

$$\text{Recall} = \frac{\#TP}{\#P} = \frac{2}{5} = 40\% \quad \text{bad!}$$

Holistic measures:

$$\text{err} = \frac{FP + FN}{n} = \frac{2 + 3}{9} = 33\% \quad \text{bad!}$$

$$F_1 = \frac{2}{\frac{1}{\text{Recall}} + \frac{1}{\text{Precision}}} = \frac{2}{\frac{1}{0.4} + \frac{1}{1}} = .57 \quad \text{bad!}$$

More measures:

False discovery rate  $FDR = \frac{FP}{\#PP}$

if  $\hat{y} = 1$ , what's your error rate?

False omission rate  $FOR = \frac{FN}{\#PN}$

if  $\hat{y} = 0$  what's your error rate

that I find the most useful and most underused!

All these can be used in-sample & OOS.

Now... sometimes each error has different costs. Example?

$$C = C_{FP} FP + C_{FN} FN$$

Idea:

Instead of err, use C as the objective function when <sup>splitting &</sup> choosing

Another option: stratified sampling: over or under sample the classes. This makes it very difficult to predict minority classes.