## Lecture 20 announcements

- HW 10 is due Thursday

- HW 6 and HW 7 grading are completed. Scores will show on D2L shortly.

## Lecture 20 outline

- CART (part 2)

- Variance of an average (e.g., an average of trees)

- Random Forest (part 1)

## CART (part2)

FOR CLASSIFICATION, USE A DIFFERENT COST FCN, E.G. :

$$\text{cost} \left\{ (\underline{x}_i, y_i) \in R_{m'} \right\}$$

$$= \frac{1}{N_{R_{m'}}} \sum_{\underline{x}_i \in R_{m'}} \mathbb{I} \left[ y_i \neq \hat{y}(R_{m'}) \right]$$

↑
CLASS
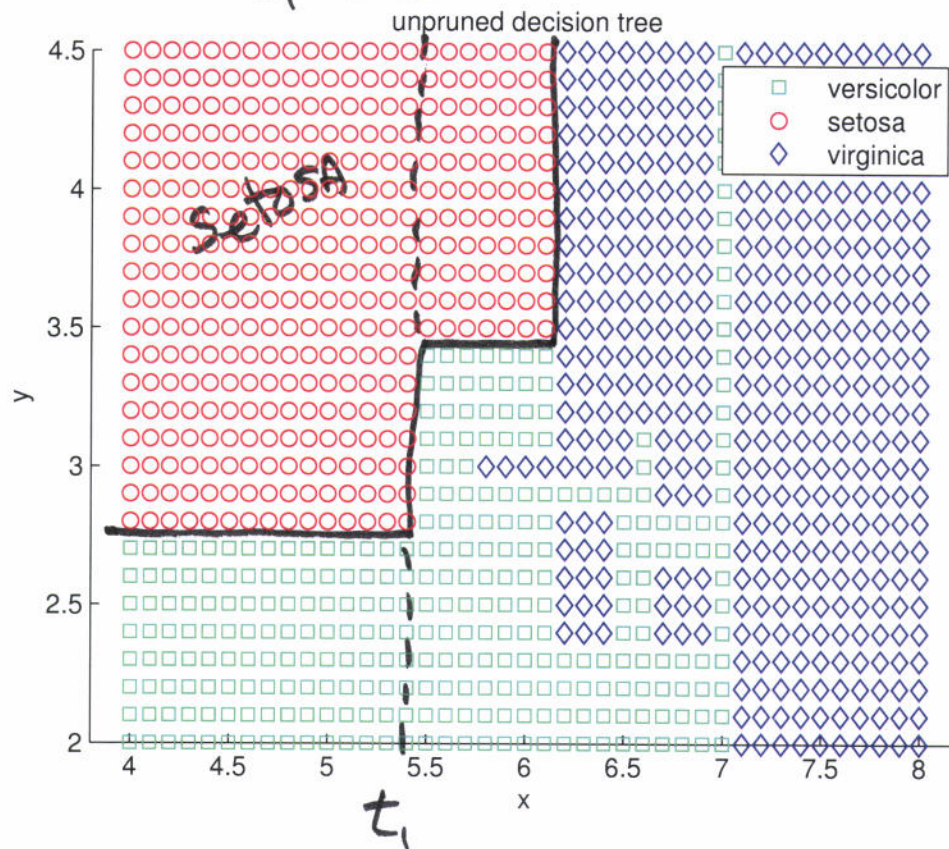ASSIGNMENT
IN $R_{m'}$

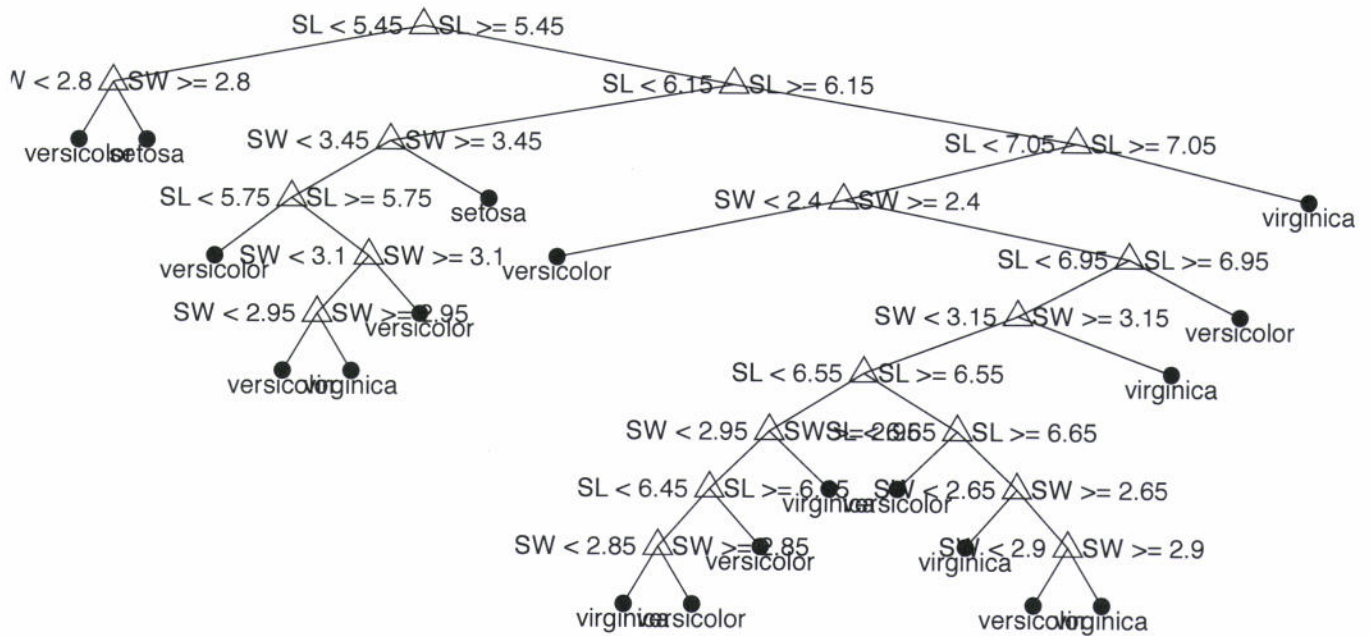(CLASS W/ MOST DATA PTS. IN $R_{m'}$)

$$= \text{CLASSIFICATION ERROR RATE} \atop \text{IN } R_{m'}.$$

OR OTHER METRIC (IN TEXT).

CART WITH THESE COST FCNS. WORKS FOR $C > 2$
CLASSES ALSO.

C = 3 CLASSES



$t = 2.8$

Murphy Fig. 16.4
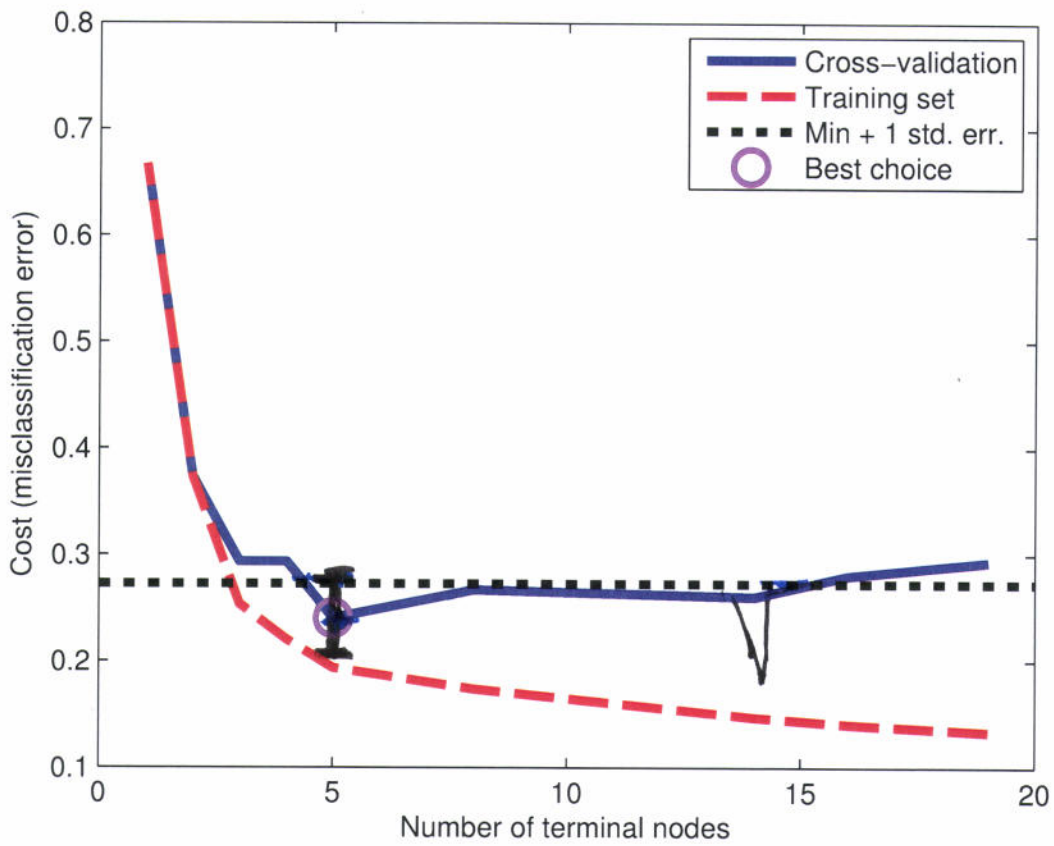
$t_1 = 5.45$

unpruned decision tree



setosa

$t_1$

Murphy Fig. 16.5(a)



Murphy Fig. 16.5(b)

[ MURPHY FIG. 16.4 - 16.5a ]
[  "     "     16.5b ]

BECAUSE CART IS A GREEDY ALGORITHM, GROWING THE
TREE UNTIL THE OPTIMAL STOPPING POINT TYPICALLY
DOESN'T YIELD THE BEST RESULTS. USUALLY IT IS
RUN PAST THIS POINT, TO YIELD A TREE THAT
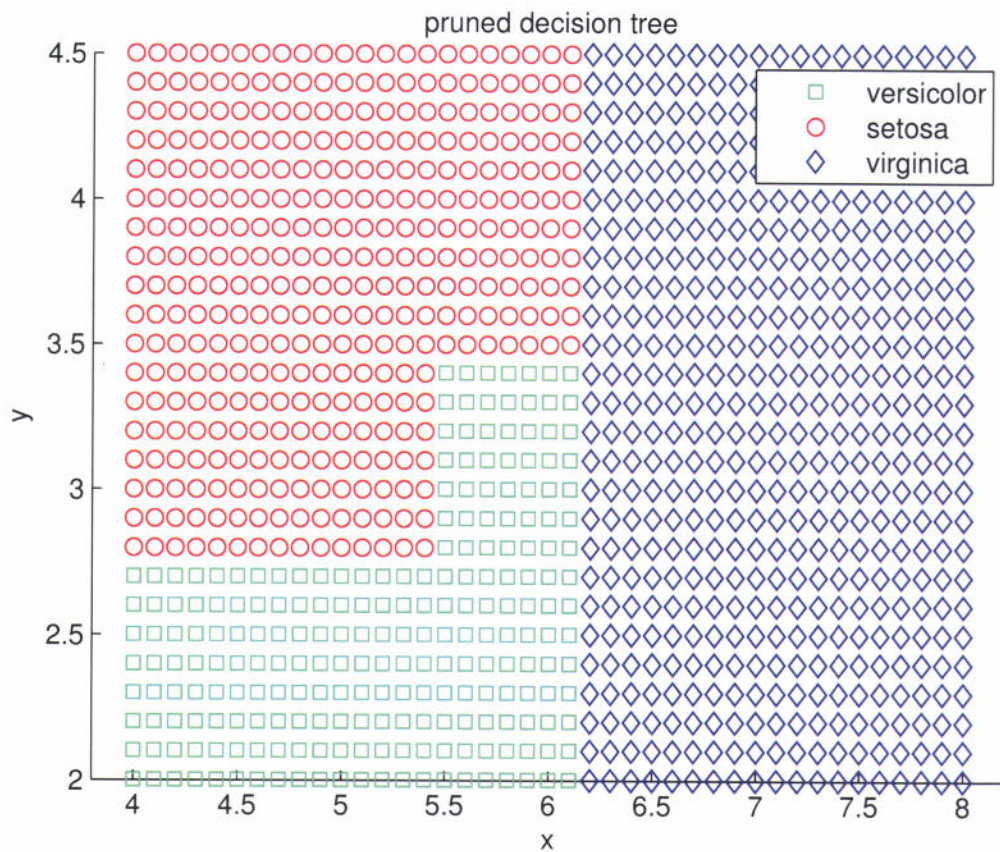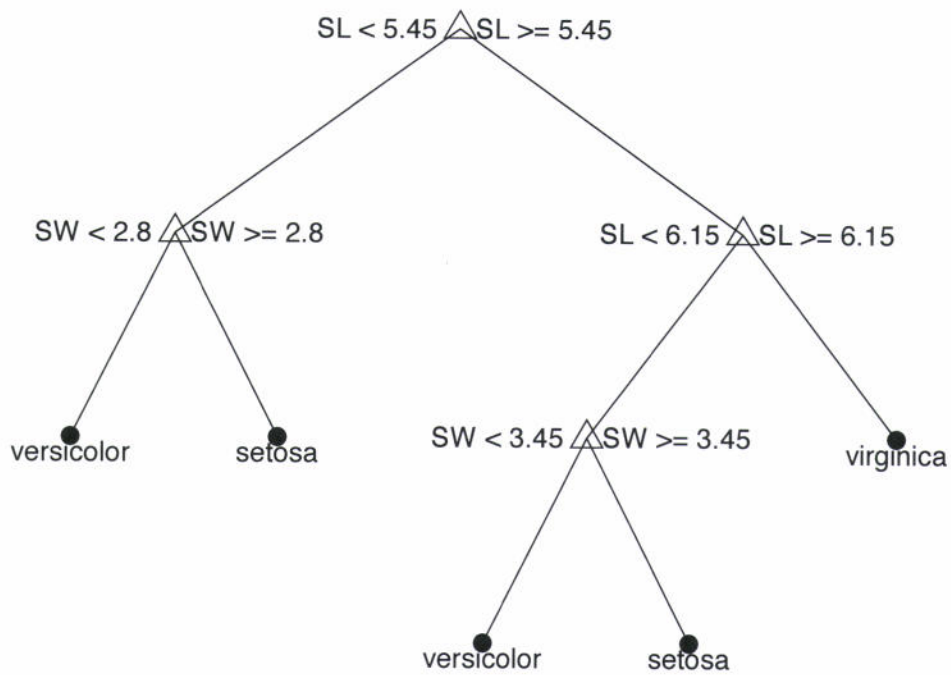OVERFITS. THEN TREE IS PRUNED.
   "WEAKEST LINK PRUNING":
     — COLLAPSE THE INTERNAL NODE THAT

     GIVES THE SMALLEST INCREASE IN
     COST FCN.; ITERATE.
     — USE CROSS-VALIDATION TO HALT WHEN
     MIN. VALIDATION ERROR IS REACHED.
     (WITHIN $1\sigma$).

     [ MURPHY FIG. 16.5b, 16.6 ].

pruned decision tree

Murphy Fig. 16.6

# CART Summary       [M 16.2.4]

| Pros | Cons |
|------|------|
| • Algorithm is fairly simple | • Predictive accuracy often isn't as good as with some other models. |
| • Incorporates linear and nonlinear boundaries on its own. | • Can be unstable to slight changes in data. |
| • Can include some automatic feature selection, | |

or Ranil feature importance.
• Typically robust to outliers.

$\Rightarrow$ HIGH VARIANCE.

## VARIANCE OF AN AVERAGE

(1)
$$\text{var}(\underline{x}) = E_{\mathscr{D}}\left\{\left(h_g^{(\mathscr{D})}(\underline{x}) - \overline{h}_g(\underline{x})\right)^2\right\}$$

$$\widetilde{h}_g^{(\mathscr{D})}(\underline{x}) \overset{\Delta}{=} \frac{1}{B}\sum_{b=1}^{B} h_g^{(\mathscr{D}_b)}(\underline{x})$$

IN WHICH $\mathscr{D}_b$ IS ONE DATASET DRAWN (WITH REPLACEMENT) FROM $\mathscr{D}$, AND USED TO ~~TO~~ TO TRAIN TREE $b$.

---

NOTE THAT:

IF WE TAKE AVERAGE OF $B$ iid RANDOM VARIABLES $v_i$, $i, = 1, 2, \cdots, B$, EACH WITH VARIANCE $\sigma^2$, THE AVERAGE WILL HAVE VARIANCE:

$$\sigma_{AVE}^2 = \frac{\sigma^2}{B}.$$

IF, INSTEAD, THE R.V. ARE IDENTICALLY DISTRIBUTED BUT HAVE PAIRWISE CORRELATION $\rho \geq 0$, ONE CAN SHOW THAT:

$$\sigma_{AVE}^2 = \rho \sigma^2 + (1-\rho) \frac{\sigma^2}{B}, \qquad 0 \le \rho \le 1$$

IN WHICH $\rho$ = CORRELATION COEFFICIENT :

(2)

$$\rho \overset{\Delta}{=} \frac{E\{v_i v_j\} - E\{v_i\} E\{v_j\}}{\sigma_{v_i} \sigma_{v_j}}$$

FOR US:

$$\rho = \frac{E\{v_i v_j\} - \mu_v^2}{\sigma_v^2}$$

$\Rightarrow$ WANT $h_g$ ~~(2)~~ $(2_6)$ TO BE AS UNCORRELATED

(BY $\rho$ MEASURE) AS POSSIBLE, AND LARGE

B.

# Random Forests

(a) Draw many datasets from $\mathcal{D}_{Tr}$, with replacement.

→ Each gives rise to a tree $T_m$ and est. $\hat{f}_m(\underline{x})$.

At each point, $\underline{x}$,

Could take average result (regression):

$$\hat{f}(x) = \frac{1}{M} \sum_{m=1}^{M} \hat{f}_m(\underline{x})$$

Or could take a vote (classification).

→ This by itself is called **BAGGING**

(for "bootstrap aggregating").

But:

- Datasets are (highly) correlated
- Reduces var some, but not a lot.

(b) BEFORE SPLITTING EACH REGION $R_m$,
 SELECT A RANDOM SUBSET OF $d$ FEATURES
 $(d < D$ OR $d << D)$,
 THEN SELECT BEST FEATURE OUT OF THE
 SUBSET TO THRESHOLD,

  $\Rightarrow$ CORRELATION BETWEEN TREES IS
   TYPICALLY MUCH SMALLER

  $\Rightarrow$ REDUCES VARIANCE BY A LOT MORE.