

Lee 19 March 2010. f 9/23/10

1

- ① Beginning with all training data....
- ② For every possible split at the current node divide the

data into X_L, \bar{y}_L & X_R, \bar{y}_R and calculate

$$SSE_L = \sum (y_L - \bar{y}_L)^2, \quad SSE_R = \sum (y_R - \bar{y}_R)^2$$

over # of
data pts under the split

- ③ Find the split with the lowest Total SSE

$$SSE_{\text{tot}} = \frac{SSE_L + SSE_R}{n_L + n_R}$$

(sorry, last time
made mistake)

- ④ Create the split, split data into two nodes
- ⑤ Repeat steps 2-4 until "STOP"

⑥ For all
leaf nodes,
assign
 $\hat{y} = \bar{y}_i$ where
 \bar{y}_i is the avg of
the y 's in that
node.

STOP: node has $\leq N_0$ data pts inside.
Default: $N_0 = 5$

There are many, many variations of the above.

If $N_0 = 1$... tree is grown to fit a separate parameter for
each data pt. $R^2 = ?$ 100%! Then tree is "pruned"

How is N_0 picked? Via an Model selection procedure

back to not
overfit.
Kind of like
backwards stepwise
regression.

train
calc lots of N_0 R^2 's or SE 's
validate....

[PBM]

Classification!

Parameter $y = \{0, 1\}$?

Binary classification

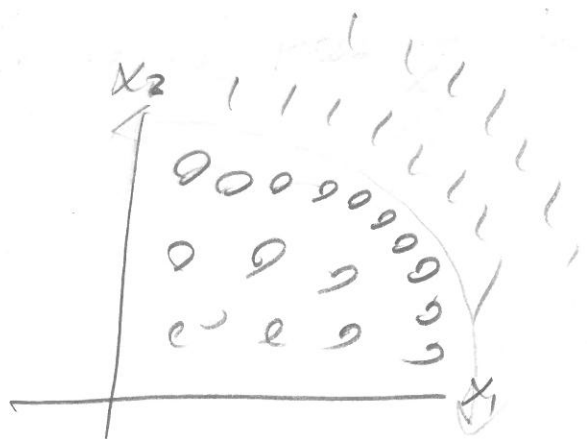
or $y \in \{1, 2, \dots, K\}$

general " into K groups / labels.

What methods do we have? For Binary classification...

- perceptron
- SVM with hinge loss

} Linear models \rightarrow could we use non-linear, polynomials, sigmoid functions? Yes



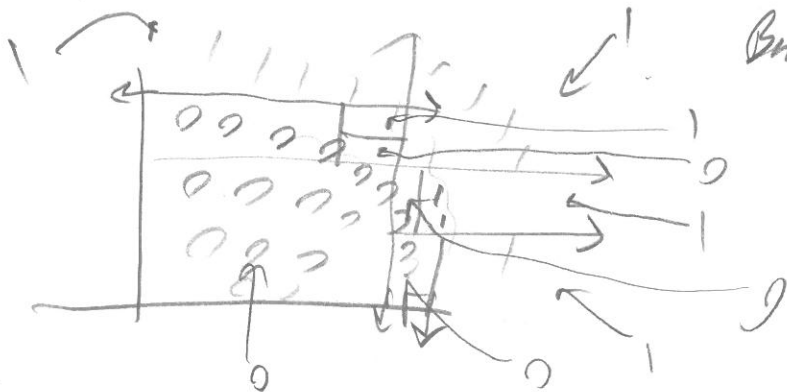
$$\mathcal{H} = \left\{ \mathbb{I}_{w_0 + w_1 x_1 + w_2 x_2 + w_3 x_1^2 + w_4 x_2^2 \geq 0} \right\}$$

$$A: \argmin_{\vec{w}} \left\{ \frac{1}{A} \sum_{i=1}^A \max \left\{ 0, \frac{1}{2} - (y_i - \frac{1}{2}) (\vec{w} \cdot \vec{x}_i - b) \right\} + \lambda \|\vec{w}\|^2 \right\}$$

Did we ever discuss general classification into $K > 2$ labels? No.

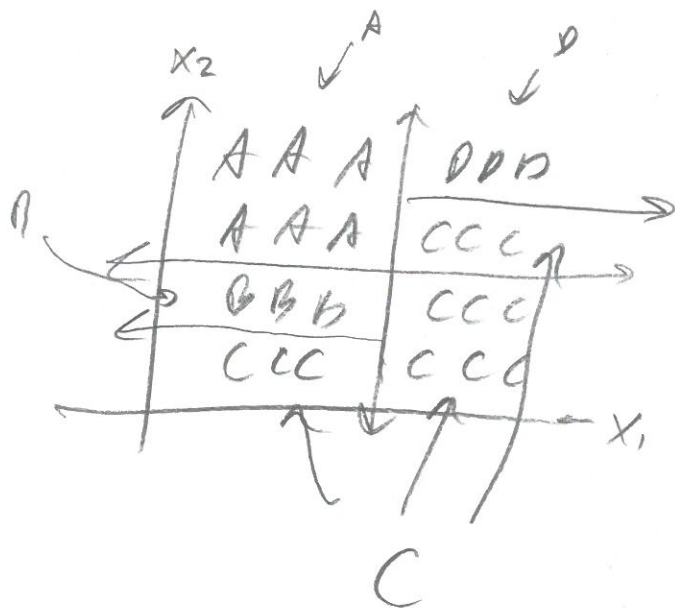
How do trees work?

Basically the same!!



Can it do more than $K=2$? Yes... very powerful..

Perceptron / SVM can do it
(as of now).



Classification Tree Algorithm

① Begin with all training data

② For every possible split, calculate the "Gini impurity" measure:

$$Gini_L := \sum_{l=1}^K \hat{p}_l(1-\hat{p}_l) \quad , \quad Gini_R := \sum_{l=1}^K \hat{p}_l(1-\hat{p}_l)$$

where $\hat{p}_l = \frac{\# y_i \text{ in category } l}{n, \# \text{ in node}}$

③ Find the split with the lowest weighted avg of Gini impurity measures:

$$Gini_{avg} := \frac{n_L Gini_L + n_R Gini_R}{n_L + n_R}$$

④ Create the split and split the data into the two categories using the two new daughter nodes

⑤ Repeat steps 2-4 until "stop" (node has $\leq N_0$ observations) Default: $N_0 = 1$!!!

⑥ For all leaf nodes, assign $\hat{y} = Mode(\vec{y}_0)$ where \vec{y}_0 are the avg. of the y_i 's in the leaf node.

DEMO