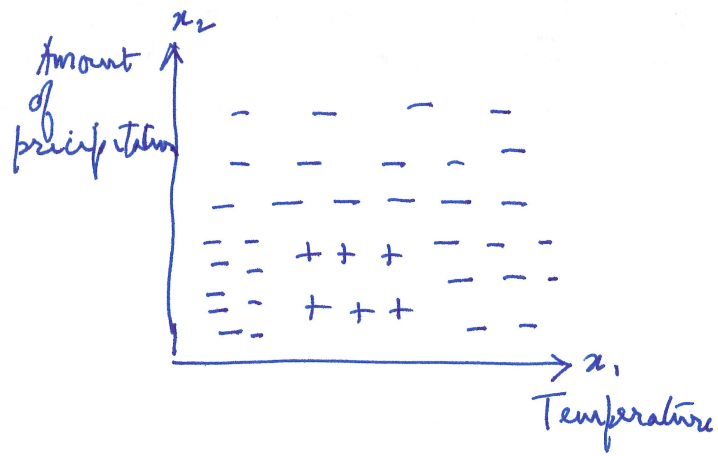
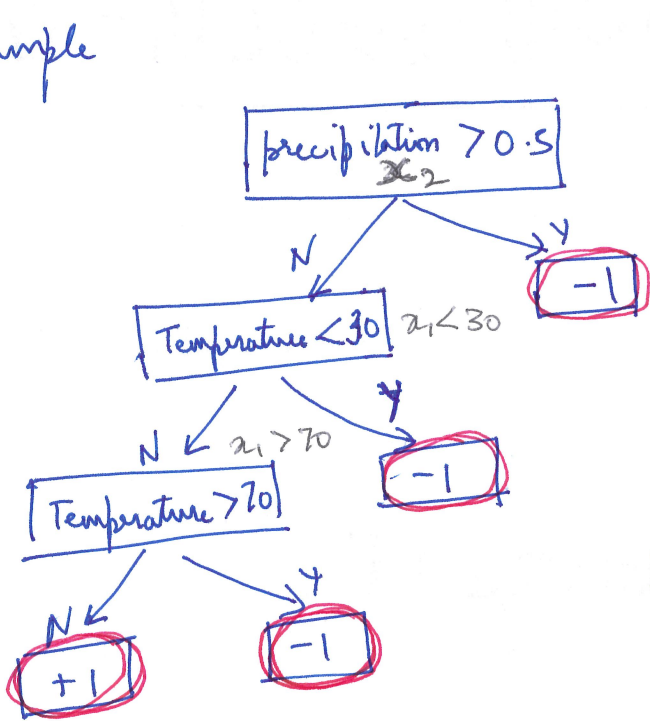


Example



$x = (50, 0.1)$  output = +1

Decision Tree: Binary tree

- Internal nodes:  $1[x_i \bowtie \theta]$  ( $\bowtie \in \{<, >, \leq, \geq\}$ )
- Leaves: Labeled by the output.

Classification: Output is  $\{+1, -1\}$

Stage k

Decision Tree  $(S, k)$

if construction is terminated ( $|S| < k$ )

Output a leaf with label = maj(S)

else

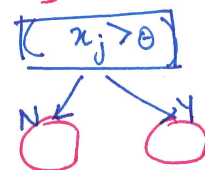
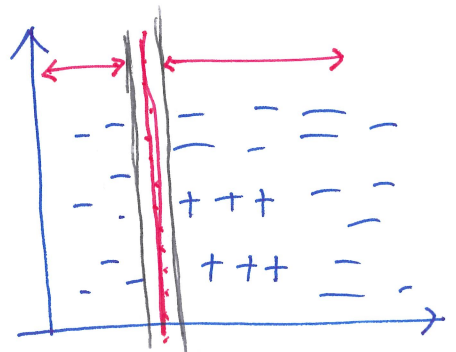
For all  $j, \theta$

$$S_N = \{x \in S \mid x_j < \theta\}$$

$$S_Y = \{x \in S \mid x_j \geq \theta\}$$

$$C(j, \theta) = (1 - \max_{a=+1, -1} P_{S_N}(a)) + (1 - \max_{a=+1, -1} P_{S_Y}(a))$$

Finite because you need to consider only finitely many values for  $\theta$ .



$P_S(a) = \frac{\# \text{ examples in } S = a}{\# S}$

pick  $j, \theta$  that minimize  $C(j, \theta)$

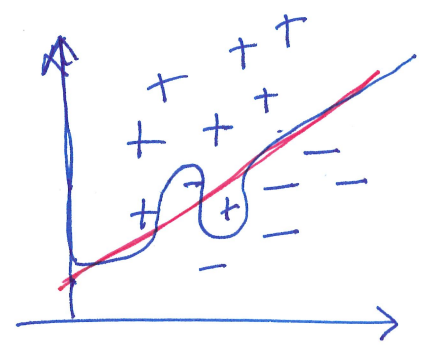
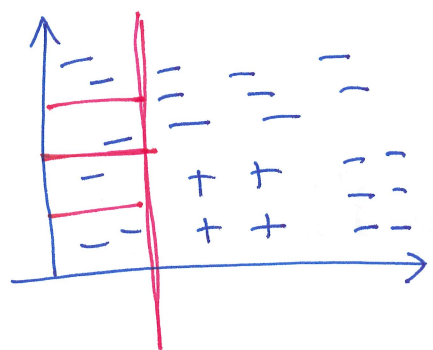
- Decision Tree  $(S_N, k)$
- Decision Tree  $(S_Y, k)$

$$\hat{p}_S(a) = \frac{\# \text{ examples in } S \text{ with output } (a)}{\# \text{ examples in } S.}$$

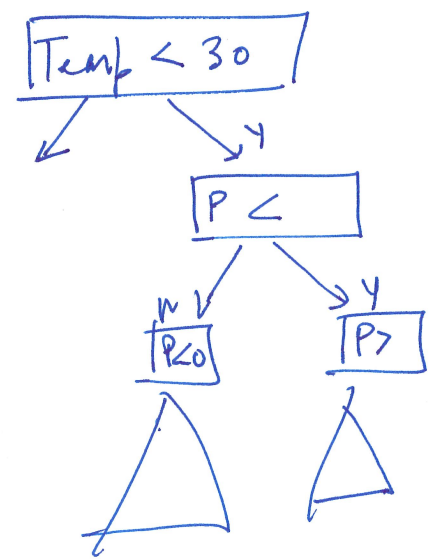
$$C_S = 1 - \max_{a \in \{1, -1\}} \hat{p}_S(a)$$

Gross entropy cost:

$$C_S = \sum_{a \in \{1, -1\}} \hat{p}_S(a) \log \hat{p}_S(a)$$



- Construction of the decision
  - Prune the decision tree
- Bottom-up process where some subtrees are replaced by leaves



Decision Trees are prone to overfitting

Bagging:

- Construct decision trees on multiple training sets
- Actual answer on a new example is "aggregation" of the answers given by each decision tree.

Construct new training sets by sampling with replacement from examples in  $S$ .

- pick some example from  $S = \{ \underbrace{(x^{(1)}, y^{(1)})}_{}, (x^{(2)}, y^{(2)}) \dots \}$