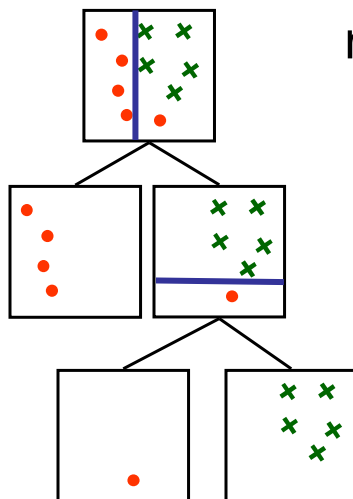


Partial Review

Decision Trees: Things to remember



1. Training: Find (variable, threshold) pairs to split the node until a “pure” node (one class) is reached
2. Test: Given new data, follow the tree until a leaf node is reached and return corresponding output value

Decision Trees

3. Measure of purity of a node →

Entropy

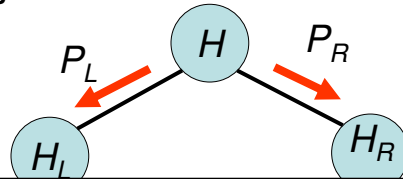
$$H(Y) = -\sum_{i=1}^n P_i \log_2 P_i$$

4. Measure of usefulness of a split →

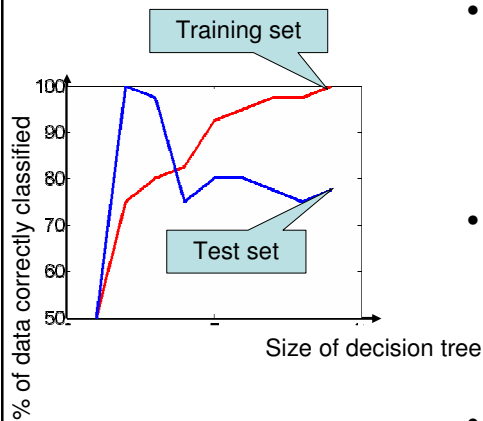
Information gain

$$IG(Y|X_j, t) = H(Y) - H(Y|X_j, t)$$

$$H(Y|X_j, t) = (H_L \times P_L + H_R \times P_R)$$



Decision Trees



- 6. Overfitting problem: DT can predict exactly Y from X from the training data but may do poorly new data
- 7. Pruning:
 - Evaluate performance on validation data
 - Remove splits if they are not statistically significant
- 8. Using pruned tree: Return the “majority” value in the leaf node reached after traversing the tree

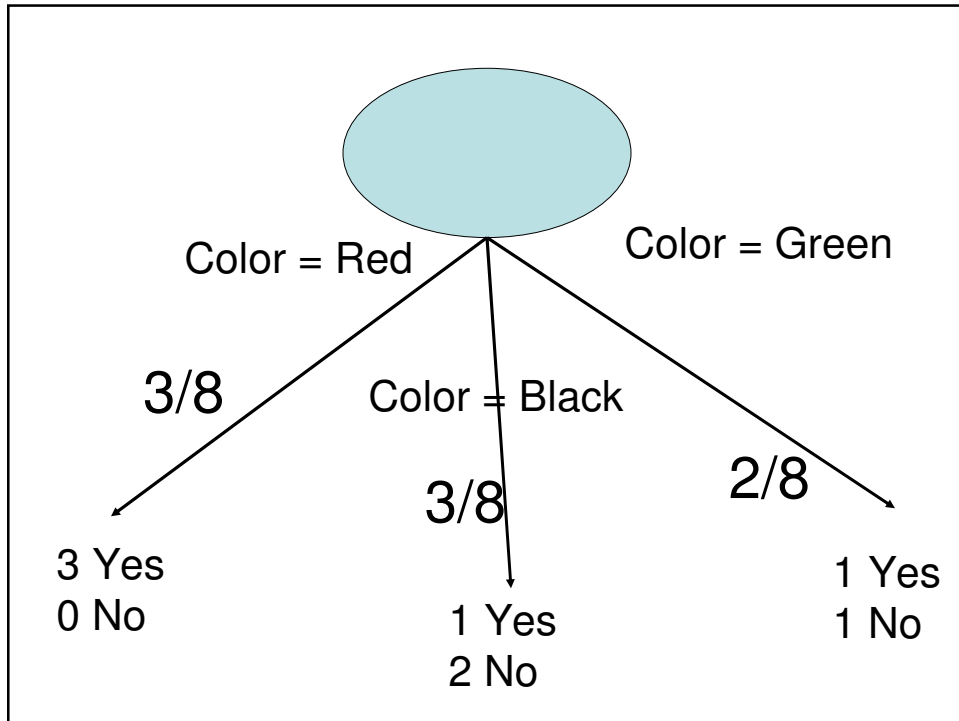
Example

Object	Color	Number of eyes	Alive
A	Red	4	Yes
B	Black	42	No
C	Red	13	Yes
D	Green	3	Yes
E	Black	27	No
F	Red	2	Yes
G	Black	1	Yes
H	Green	11	No

- $H(\text{Alive}) = ?$
- $IG(\text{Alive}|\text{Color}) = ?$

Example

- $H(\text{Alive}) = -P_{\text{Yes}} \log_2(P_{\text{Yes}}) - P_{\text{No}} \log_2(P_{\text{No}})$
- $P_{\text{Yes}} = 5/8$
- $P_{\text{No}} = 3/8$
- $H(\text{Alive}) \sim 0.95$



Example

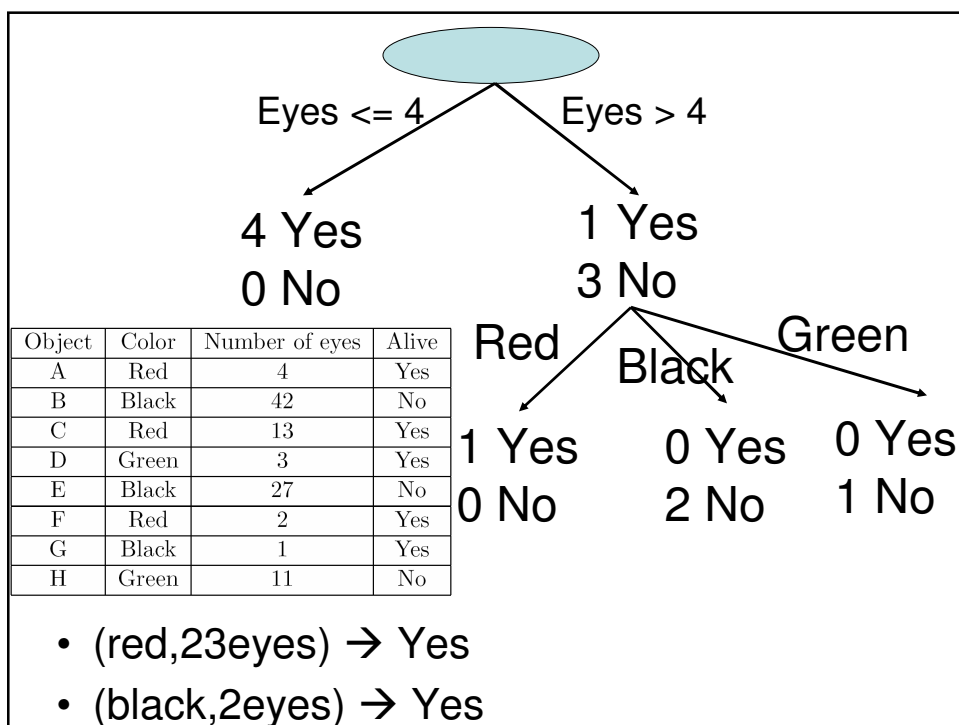
Color = Red Color = Black Color = Green

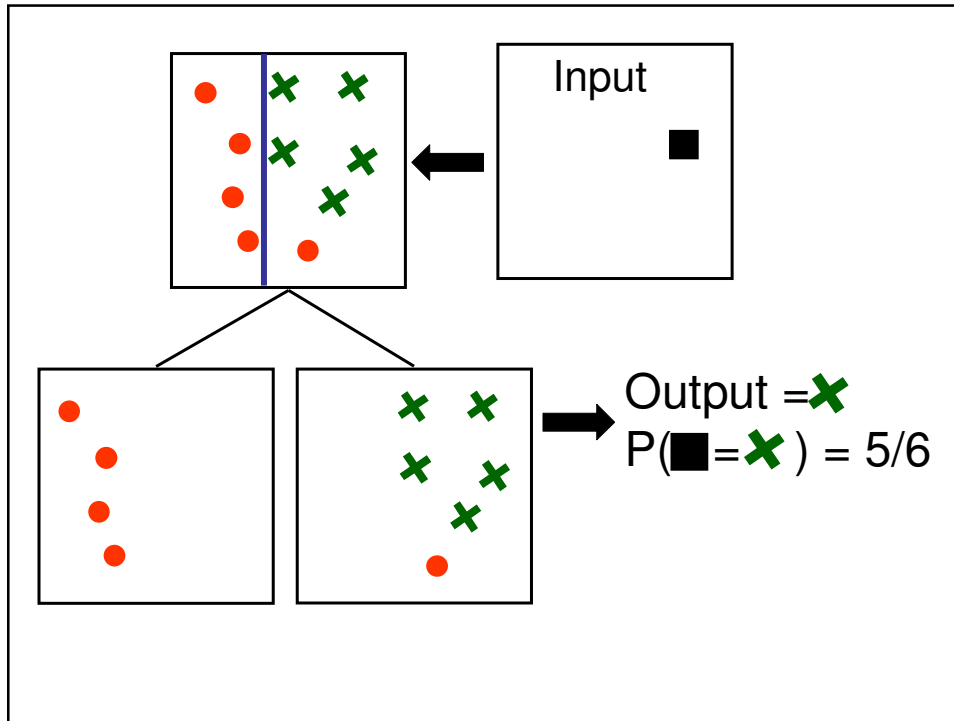
3 Yes 1 Yes 1 Yes
0 No 2 No 1 No

- $H(\text{Alive}|\text{Color} = \text{Red}) = 0$
- $H(\text{Alive}|\text{Color} = \text{Black}) = -1/3 \log(1/3) - 2/3 \log(2/3)$
- $H(\text{Alive}|\text{Color} = \text{Green}) = -1/2 \log(1/2) - 1/2 \log(1/2) = 1$
- $IG(\text{Alive}|\text{Color}) = H(\text{Alive}) - 3/8 \times 0 - 3/8 \times H(\text{Alive}|\text{Color} = \text{Black}) - 2/8 \times 1 \sim 0.36$

Example

- Tree when splitting first on (Eyes > 4) and then on color?
- Classification for data = (red,23eyes)?
- Classification for data = (black,2eyes)?





Neural Nets: Things to remember

Input attribute values

Output prediction

$\sum_{i=0}^M w_i x_i = \mathbf{w} \cdot \mathbf{x}$

σ

1. Output prediction is either $\mathbf{w} \cdot \mathbf{x}$ or $\sigma(\mathbf{w} \cdot \mathbf{x})$
2. Can solve any *linear* problem
3. Training by minimizing:

$$E = \sum_{\text{Training Data}} (\mathbf{y}^k - \mathbf{w} \cdot \mathbf{x}^k)^2 = \sum_{k=1}^N \delta_k^2$$

Neural Nets

4. Incremental update rule during training (guaranteed convergence):

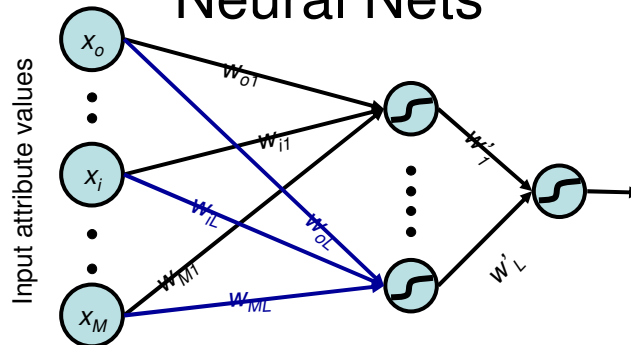
$$w_i \leftarrow w_i - \alpha \frac{\partial E}{\partial w_i} \quad \rightarrow \quad w_i \leftarrow w_i + \alpha \delta_k x_i^k$$

OR

$$w_i \leftarrow w_i + \alpha \delta_k x_i^k \sigma'(\mathbf{w} \cdot \mathbf{x}^k)$$

5. α = learning rate (low = smooth convergence but slow, high = chaotic convergence but fast)

Neural Nets



6. Can represent any function of the input attributes (given enough layers and units)
7. Convergence is *not* guaranteed

Example

- Inputs: M binary variables x_1, \dots, x_M (M odd)
- Output: $Y =$ majority value of the inputs, i.e., $Y = 1$ if there are more 1s than zeros, $Y = 0$ otherwise
- Design a neural networks that computes Y from x_1, \dots, x_M

