#### **Decision Trees**

Aarti Singh

Machine Learning 10-315 Oct 6, 2021

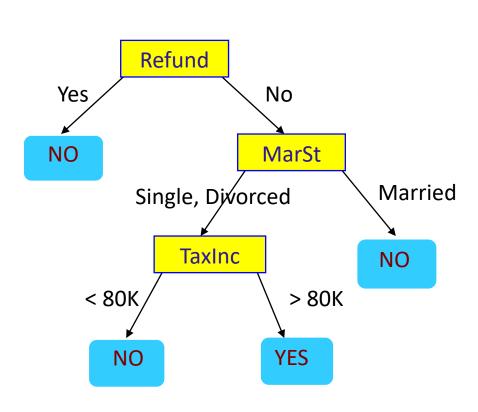


### **Decision Trees**

• Start with discrete features, then discuss continuous

## Representation

What does a decision tree represent

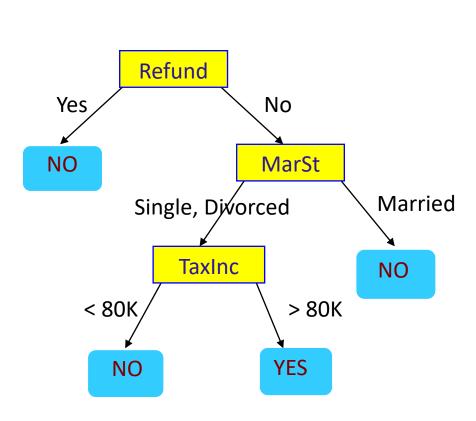


$X_1$	$X_2$	$X_3$	Y
Refund	Marital Status	Taxable Income	Cheat

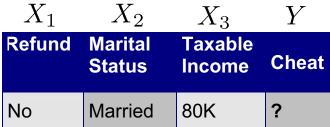
- Each internal node: test one feature X<sub>i</sub>
- Each branch from a node: selects some value for X<sub>i</sub>
- Each leaf node: prediction for Y

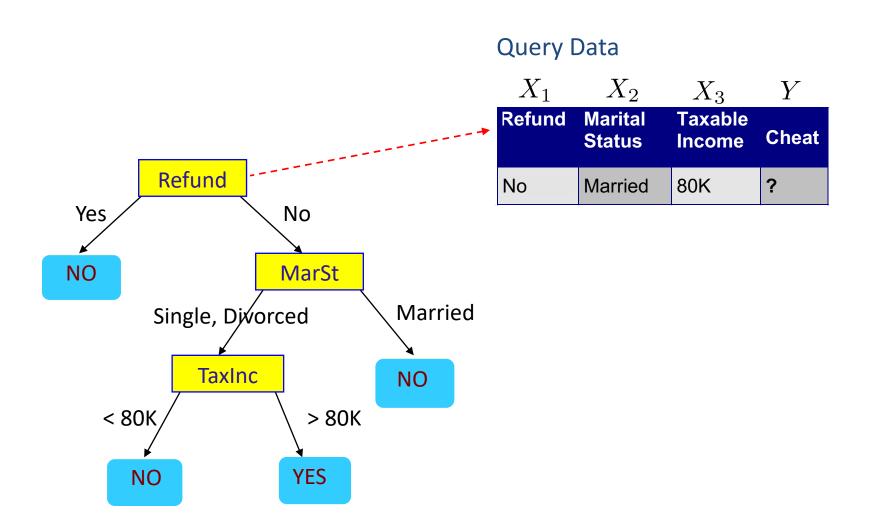
### **Prediction**

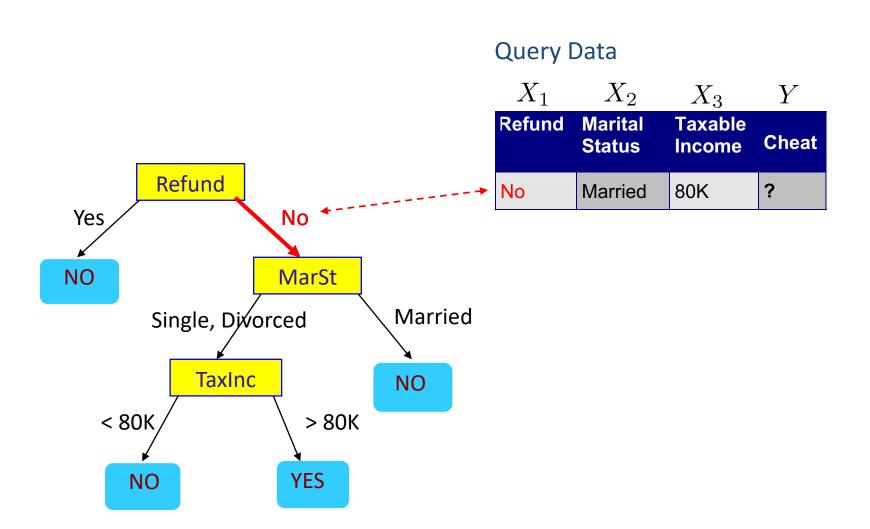
 Given a decision tree, how do we assign label to a test point

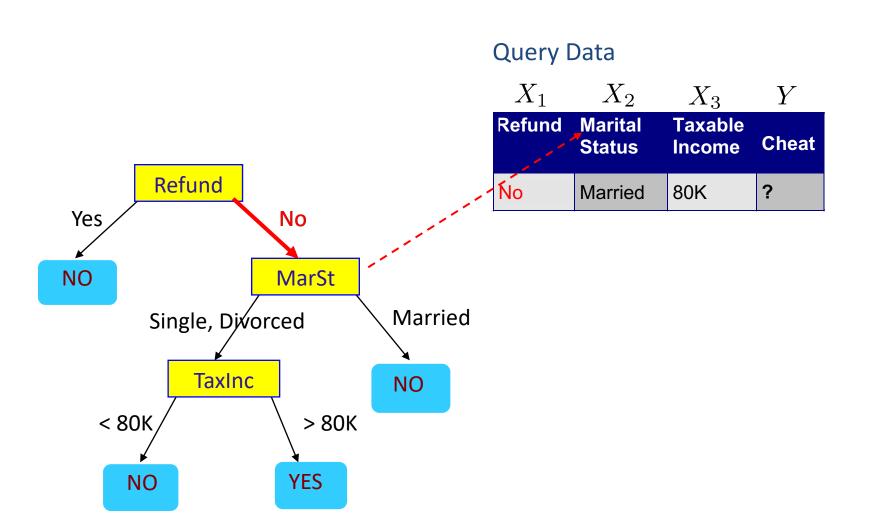


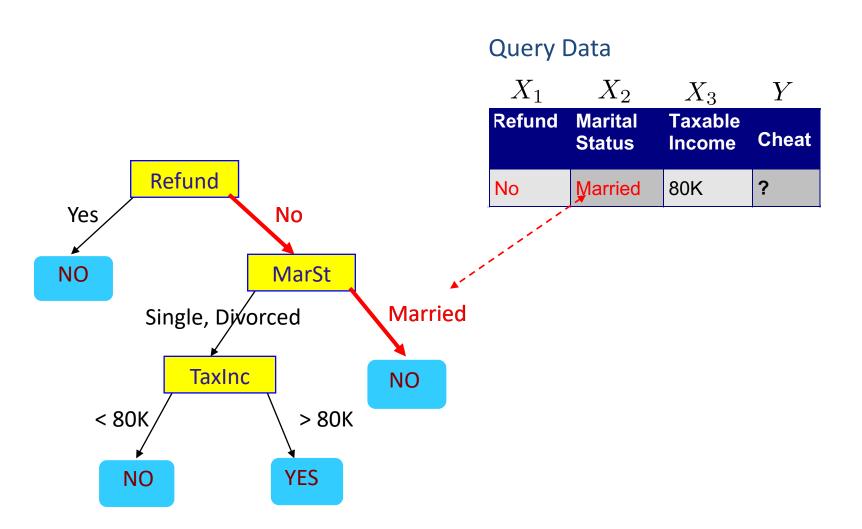
#### **Query Data**

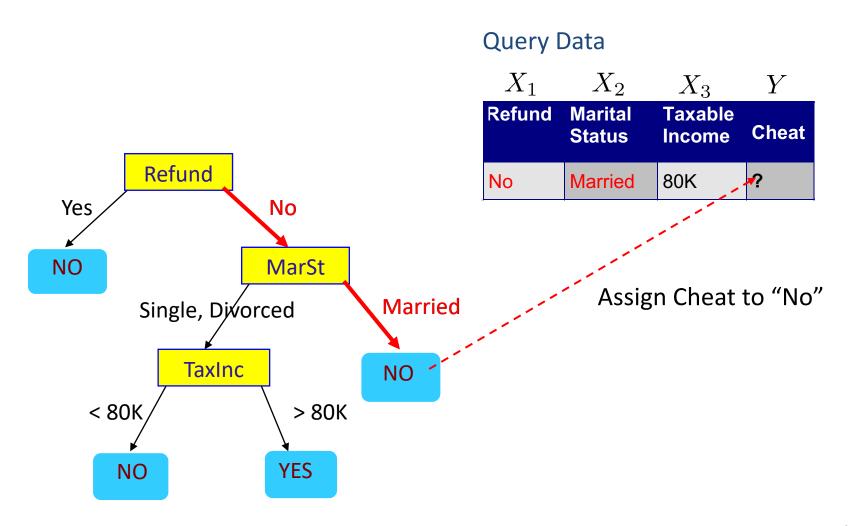












#### So far...

- What does a decision tree represent
- Given a decision tree, how do we assign label to a test point

Discriminative or Generative?

#### Now ...

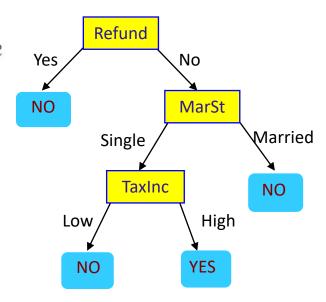
How do we learn a decision tree from training data

#### How to learn a decision tree

Top-down induction [ID3]

#### Main loop:

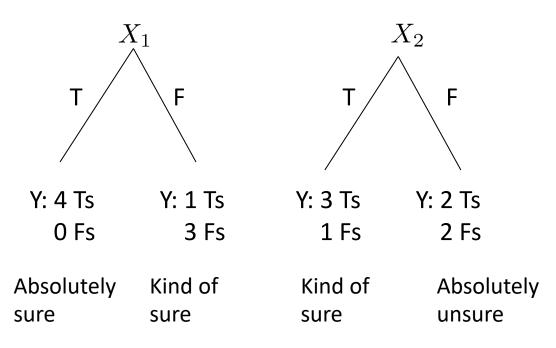
- 1.  $X \leftarrow$  the "best" decision feature—for next node
- 2. Assign X as decision feature—for node
- 3. For each value of X, create new descendant of node (Discrete features)
- 4. Sort training examples to leaf nodes
- 5. If training examples perfectly classified, Then STOP, Else iterate over new leaf nodes (steps 1-5) after removing current feature



6. When all features exhausted, assign majority label to the leaf node

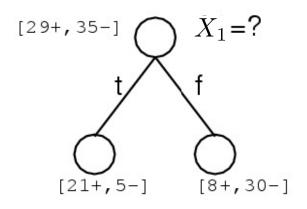
#### Which feature is best?

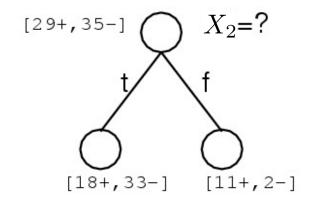
X <sub>1</sub>	$X_2$	Υ
Τ	_	Т
Τ	F	Т
Т	Т	Т
Т	F	Т
F	F T	
F F		F
F	Т	F
F	F	F



Good split if we are more certain about classification after split – Uniform distribution of labels is bad

#### Which feature is best?





Pick the attribute/feature which yields maximum information gain:

$$\arg\max_{i} I(Y, X_i) = \arg\max_{i} [H(Y) - H(Y|X_i)]$$

H(Y) – entropy of Y  $H(Y|X_i)$  – conditional entropy of Y

## **Andrew Moore's Entropy in a Nutshell**



Low Entropy

**High Entropy** 

..the values (locations of soup) sampled entirely from within the soup bowl

..the values (locations of soup) unpredictable... almost uniformly sampled throughout our dining room

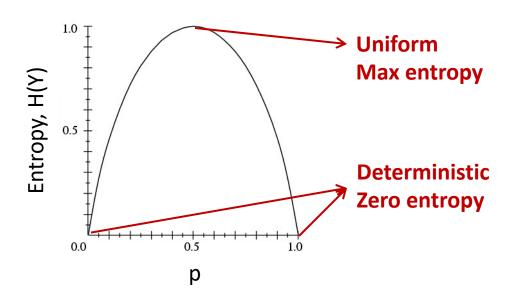
## **Entropy**

Entropy of a random variable Y

$$H(Y) = -\sum_{y} P(Y = y) \log_2 P(Y = y)$$

More uncertainty, more entropy!

Y ~ Bernoulli(p)



# **Information Theory interpretation**

**Entropy**: H(Y) = H(P) is the expected number of bits needed to encode a randomly drawn value of  $Y \sim P$  under most efficient code optimized for distribution P

**Cross-Entropy**: H(P,Q) is the expected number of bits needed to encode a randomly drawn value of Y under most efficient code optimized for distribution Q

#### **Information Gain**

- Advantage of attribute = decrease in uncertainty
  - Entropy of Y before split

$$H(Y) = -\sum_{y} P(Y = y) \log_2 P(Y = y)$$

- Entropy of Y after splitting based on X<sub>i</sub>
  - Weight by probability of following each branch

$$H(Y \mid X_i) = \sum_{x} P(X_i = x) H(Y \mid X_i = x)$$
  
=  $-\sum_{x} P(X_i = x) \sum_{y} P(Y = y \mid X_i = x) \log_2 P(Y = y \mid X_i = x)$ 

Information gain is difference

$$I(Y, X_i) = H(Y) - H(Y \mid X_i)$$

Max Information gain = min conditional entropy

# Which feature is best to split?

Pick the attribute/feature which yields maximum information gain:

$$rg \max_i I(Y,X_i) = rg \max_i [H(Y) - H(Y|X_i)]$$
 
$$= rg \min_i H(Y|X_i)$$
 Entropy of Y 
$$H(Y) = -\sum_i P(Y=y) \log_2 P(Y=y)$$

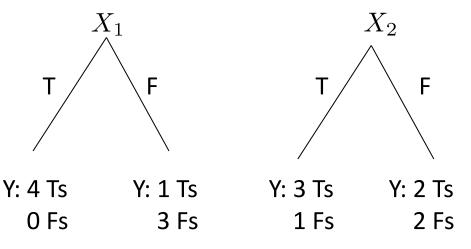
Conditional entropy of Y  $H(Y \mid X_i) = \sum_{x} P(X_i = x) H(Y \mid X_i = x)$ 

Feature which yields maximum reduction in entropy (uncertainty) provides maximum information about Y

#### **Information Gain**

$$H(Y \mid X_i) = -\sum_{x} P(X_i = x) \sum_{y} P(Y = y \mid X_i = x) \log_2 P(Y = y \mid X_i = x)$$

X <sub>1</sub>	$X_2$	Υ	
Η	T		
Η	F T		
Η	T		
Τ	F	Т	
F	Τ	Т	
F	F	F	
F	Т	T F	
F	F	F	



$$\widehat{H}(Y|X_1) = -\frac{1}{2} \left[1 \log_2 1 + 0 \log_2 0\right] - \frac{1}{2} \left[\frac{1}{4} \log_2 \frac{1}{4} + \frac{3}{4} \log_2 \frac{3}{4}\right]$$

$$\widehat{H}(Y|X_2) = -\frac{1}{2} \left[\frac{3}{4} \log_2 \frac{3}{4} + \frac{1}{4} \log_2 \frac{1}{4}\right] - \frac{1}{2} \left[\frac{1}{2} \log_2 \frac{1}{2} + \frac{1}{2} \log_2 \frac{1}{2}\right]$$

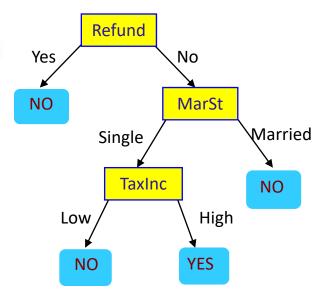
$$\widehat{H}(Y|X_1) < \widehat{H}(Y|X_2)$$

#### How to learn a decision tree

Top-down induction [ID3]

#### Main loop:

- 1.  $X \leftarrow$  the "best" decision feature—for next node
- 2. Assign X as decision feature—for node
- 3. For each value of X, create new descendant of node (Discrete features)
- 4. Sort training examples to leaf nodes
- 5. If training examples perfectly classified, Then STOP, Else iterate over new leaf nodes (steps 1-5) after removing current feature



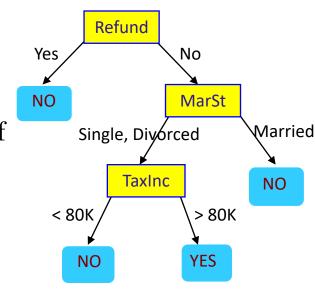
6. When all features exhausted, assign majority label to the leaf node

#### How to learn a decision tree

• Top-down induction [ID3, C4.5, C5, ...]

#### Main loop: C4.5

- 1.  $X \leftarrow$  the "best" decision feature—for next node
- 2. Assign X as decision feature—for node
- 3. For "best" split of X, create new descendants of node
- 4. Sort training examples to leaf nodes
- 5. If training examples perfectly classified, Then STOP, Else iterate over new leaf nodes
- 6. Prune back tree to reduce overfitting
- 7. Assign majority label to the leaf node



## Handling continuous features (C4.5)

Convert continuous features into discrete by setting a threshold.

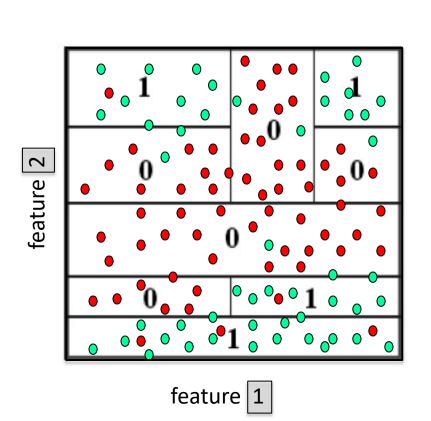
What threshold to pick?

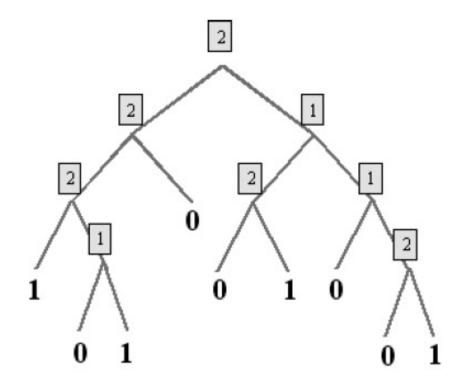
Search for best one as per information gain. Infinitely many??

Don't need to search over more than  $\sim$  n (number of training data),e.g. say  $X_1$  takes values  $x_1^{(1)}$ ,  $x_1^{(2)}$ , ...,  $x_1^{(n)}$  in the training set. Then possible thresholds are

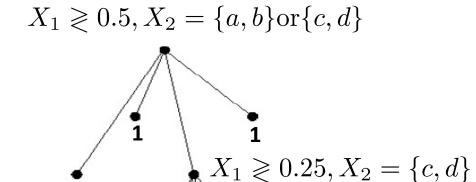
$$[x_1^{(1)} + x_1^{(2)}]/2$$
,  $[x_1^{(2)} + x_1^{(3)}]/2$ , ...,  $[x_1^{(n-1)} + x_1^{(n)}]/2$ 

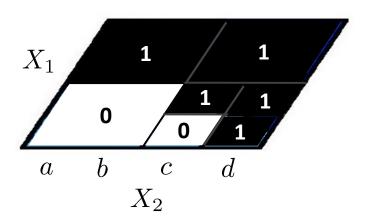
# Dyadic decision trees (split on mid-points of features)





## **Decision Tree more generally...**

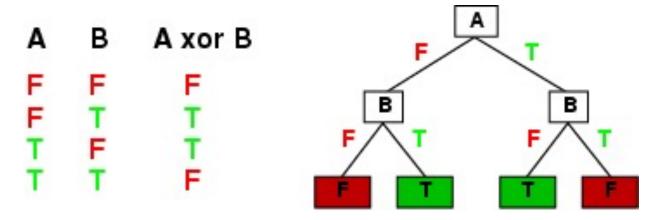




- Features can be discrete, continuous or categorical
- Each internal node: test some set of features {X<sub>i</sub>}
- Each branch from a node: selects a set of value for {X<sub>i</sub>}
- Each leaf node: prediction for Y

## **Expressiveness of Decision Trees**

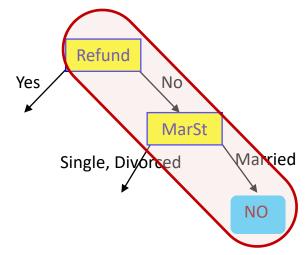
- Decision trees in general (without pruning) can express any function of the input features.
- E.g., for Boolean functions, truth table row → path to leaf:



- There is a decision tree which perfectly classifies a training set with one path to leaf for each example - overfitting
- But it won't generalize well to new examples prefer to find more compact decision trees

## When to Stop?

- Many strategies for picking simpler trees:
  - Pre-pruning
    - Fixed depth (e.g. ID3)
    - Fixed number of leaves
  - Post-pruning
    - Chi-square test
      - Convert decision tree to a set of rules
      - Eliminate variable values in rules which are independent of label (using chi-square test for independence)
      - Simplify rule set by eliminating unnecessary rules
  - Information Criteria: MDL(Minimum Description Length)



## **Information Criteria**

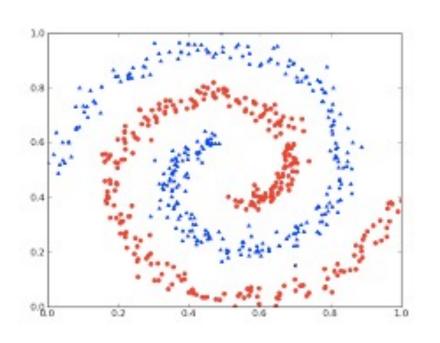
Penalize complex models by introducing cost

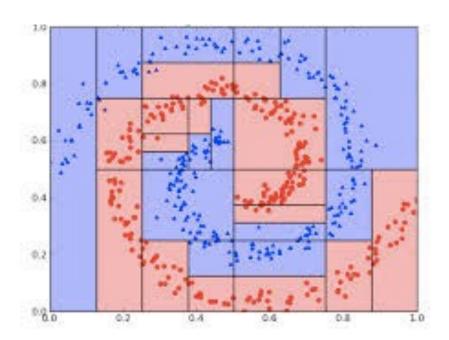
$$\widehat{f} = \arg\min_{T} \ \left\{ \frac{1}{n} \sum_{i=1}^{n} \mathsf{loss}(\widehat{f}_{T}(X_{i}), Y_{i}) \ + \ \mathsf{pen}(T) \right\}$$
 
$$\mathsf{log} \ \mathsf{likelihood} \qquad \mathsf{cost}$$

$$loss(\widehat{f}_T(X_i), Y_i) = (\widehat{f}_T(X_i) - Y_i)^2$$
 regression  $= \mathbf{1}_{\widehat{f}_T(X_i) \neq Y_i}$  classification

 $pen(T) \propto |T|$  penalize trees with more leaves CART – optimization can be solved by dynamic programming

# Example of 2-feature decision tree classifier

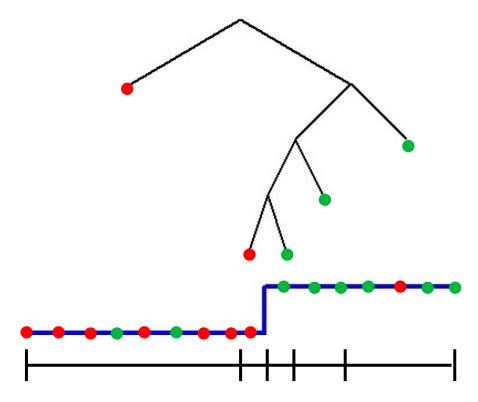




## How to assign label to each leaf

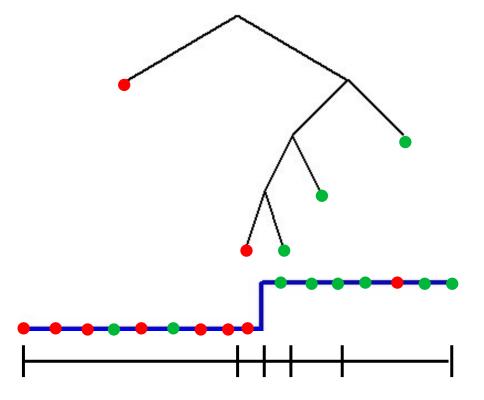
Classification – Majority vote

Regression –?

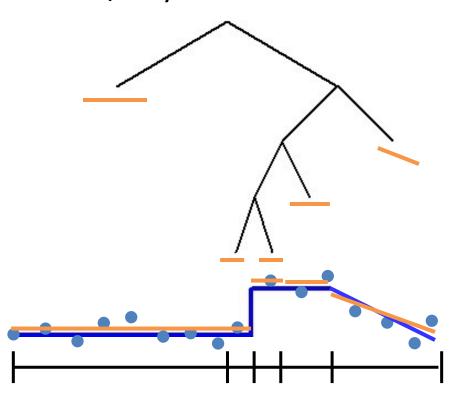


# How to assign label to each leaf

Classification – Majority vote



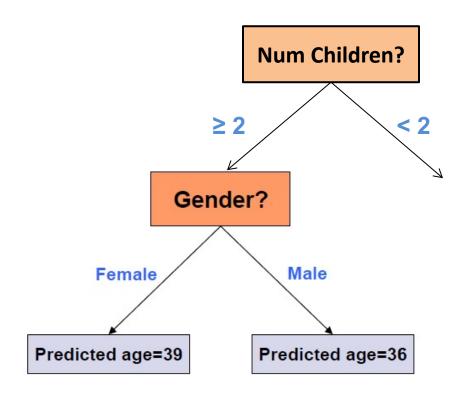
Regression – Constant/ Linear/Poly fit



## Regression trees

 $X^{(1)}$  ....  $X^{(p)}$  Y

Gender	Rich?	Num. Children	# travel per yr.	Age
F	No	2	5	38
М	No	0	2	25
М	Yes	1	0	72
:	:		:	:



Average (fit a constant ) using training data at the leaves

## What you should know

- Decision trees are one of the most popular data mining tools
  - Simplicity of design
  - Interpretability
  - Ease of implementation
  - Good performance in practice (for small dimensions)
- Information gain to select attributes (ID3, C4.5,...)
- Decision trees will overfit!!!
  - Must use tricks to find "simple trees", e.g.,
    - Pre-Pruning: Fixed depth/Fixed number of leaves
    - Post-Pruning: Chi-square test of independence
    - Complexity Penalized/MDL model selection
- Can be used for classification, regression and density estimation too