



# Big Data Analysis: Classification - Decision Tree

---

Younghoon Kim

CS246: Mining Massive Datasets  
Jure Leskovec, Stanford University  
<http://cs246.stanford.edu>

# Decision Trees on MapReduce

CS246: Mining Massive Datasets  
Jure Leskovec, Stanford University  
<http://cs246.stanford.edu>





# Classification

---

## ■ Given

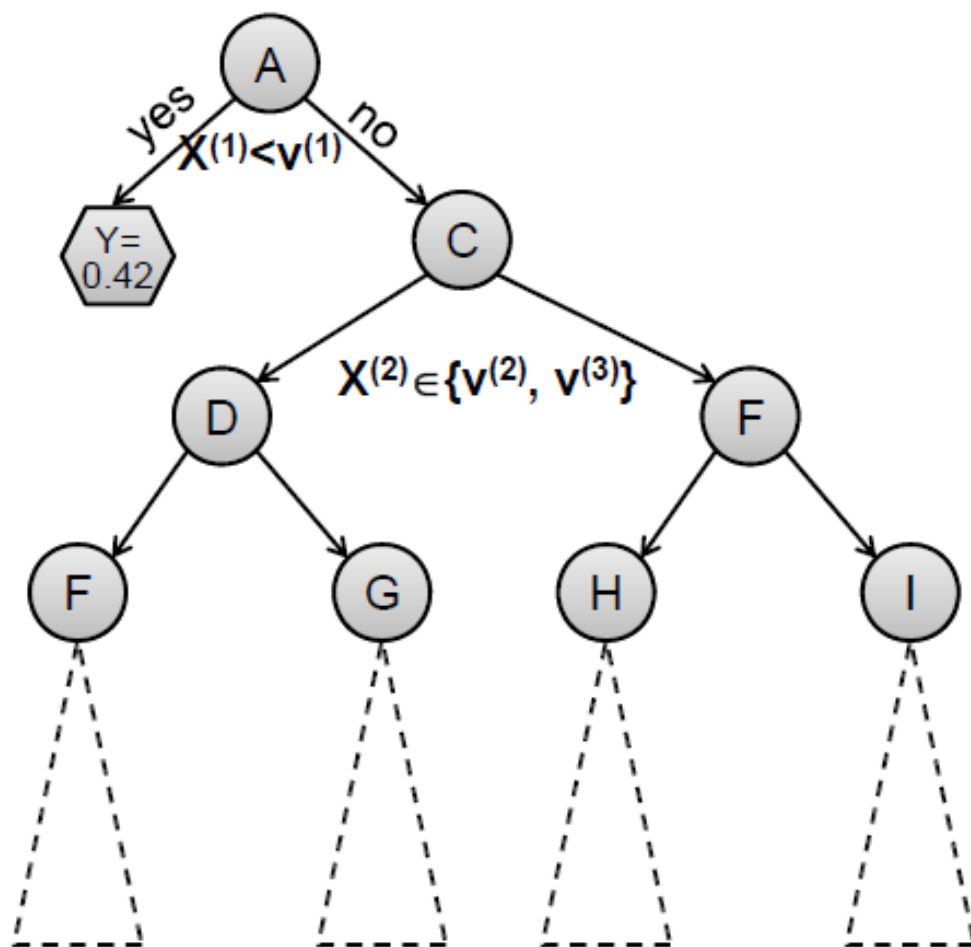
- A set of  $d$ -dimensional vectors:  $x_i = \langle x_i^{(1)}, x_i^{(2)}, \dots, x_i^{(d)} \rangle$
- $O_j$ : domain of  $x^{(j)}$ 
  - Categorical: e.g.,  $O_j = \{\text{red, blue, orange, green}\}$
  - Numerical:  $O_j = \{1, 2, \dots, 10\}$  or  $\mathbb{R}$
- $y_i$ : a label in a categorical domain  $O_Y$ 
  - E.g.,  $\{\text{yes, no}\}$  or  $\{+, -\}$
- Data  $D: \{x_i, y_i\}_{i=1, \dots, n}$

## ■ Goal

- Given a  $d$ -dimensional vector  $x$ , whose  $y$  is unknown,
- Predict  $y$

# Decision Trees

- A **Decision Tree** is a tree-structured plan of a set of attributes to test in order to predict the output



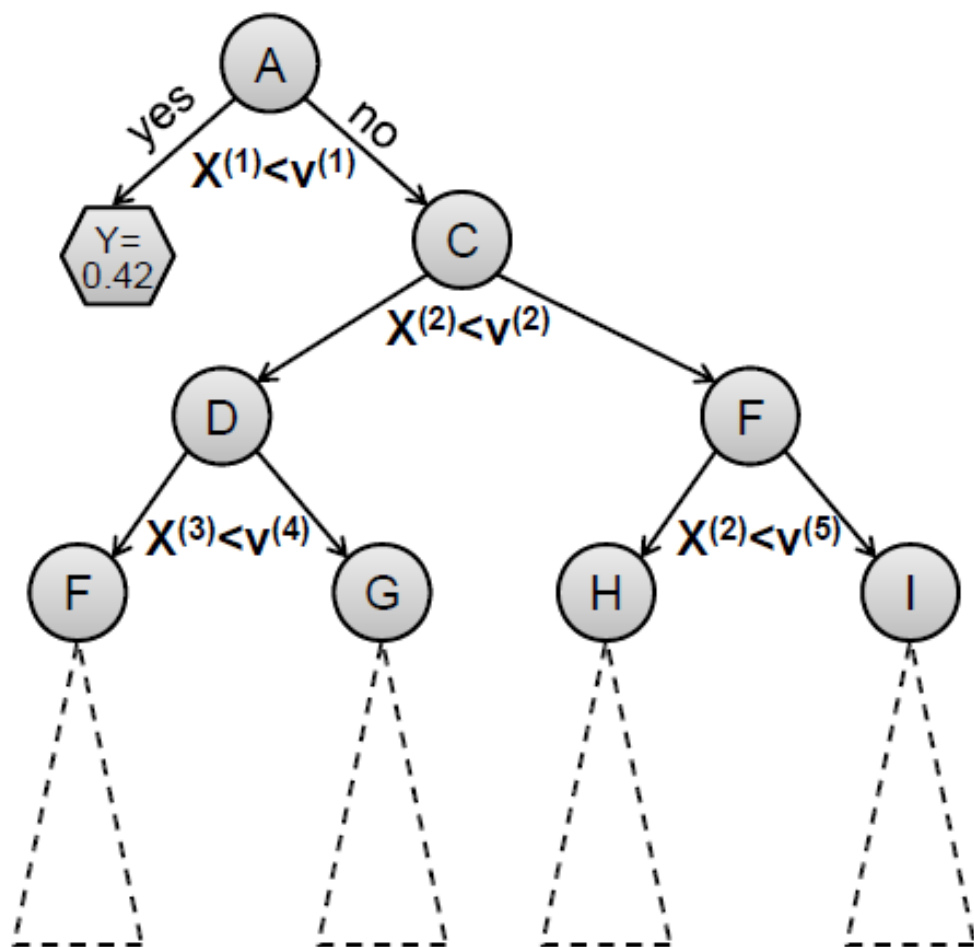
# Decision Trees (1)

## ■ Decision trees:

- Split the data at each internal node
- Each leaf node makes a prediction

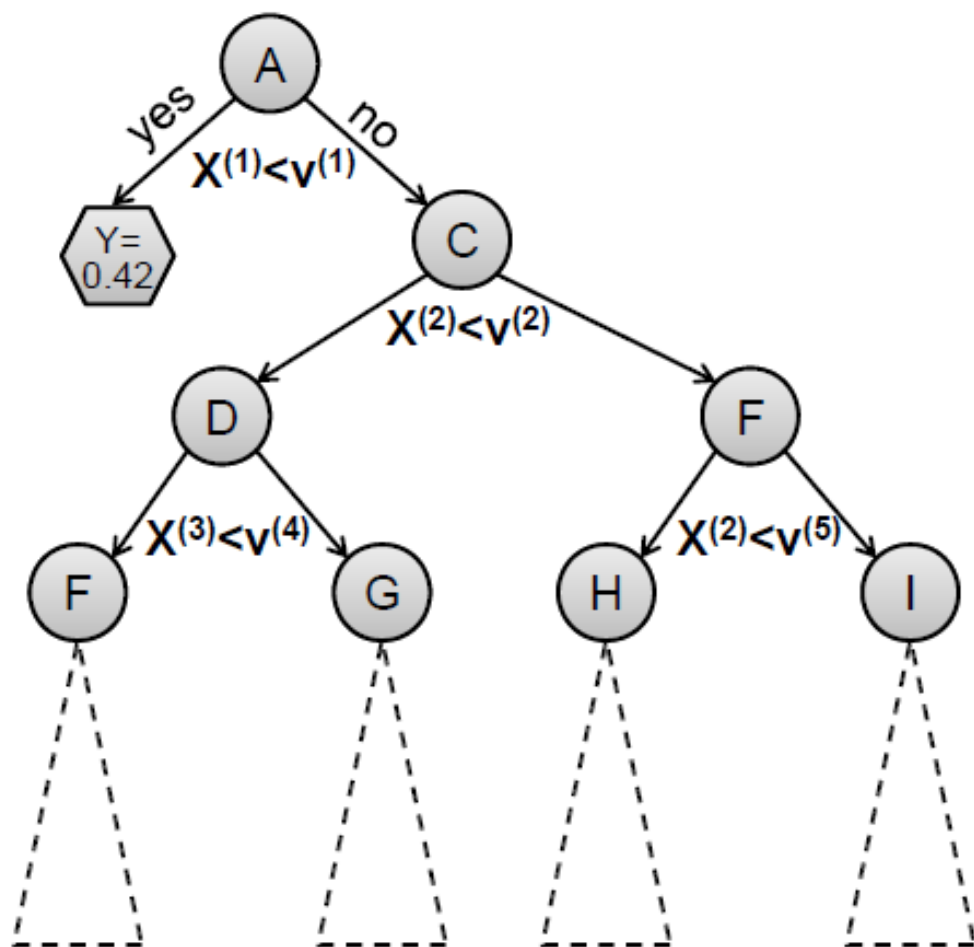
## ■ Lecture today:

- Binary splits:  $X^{(j)} < v$
- Numerical attrs.
- Regression



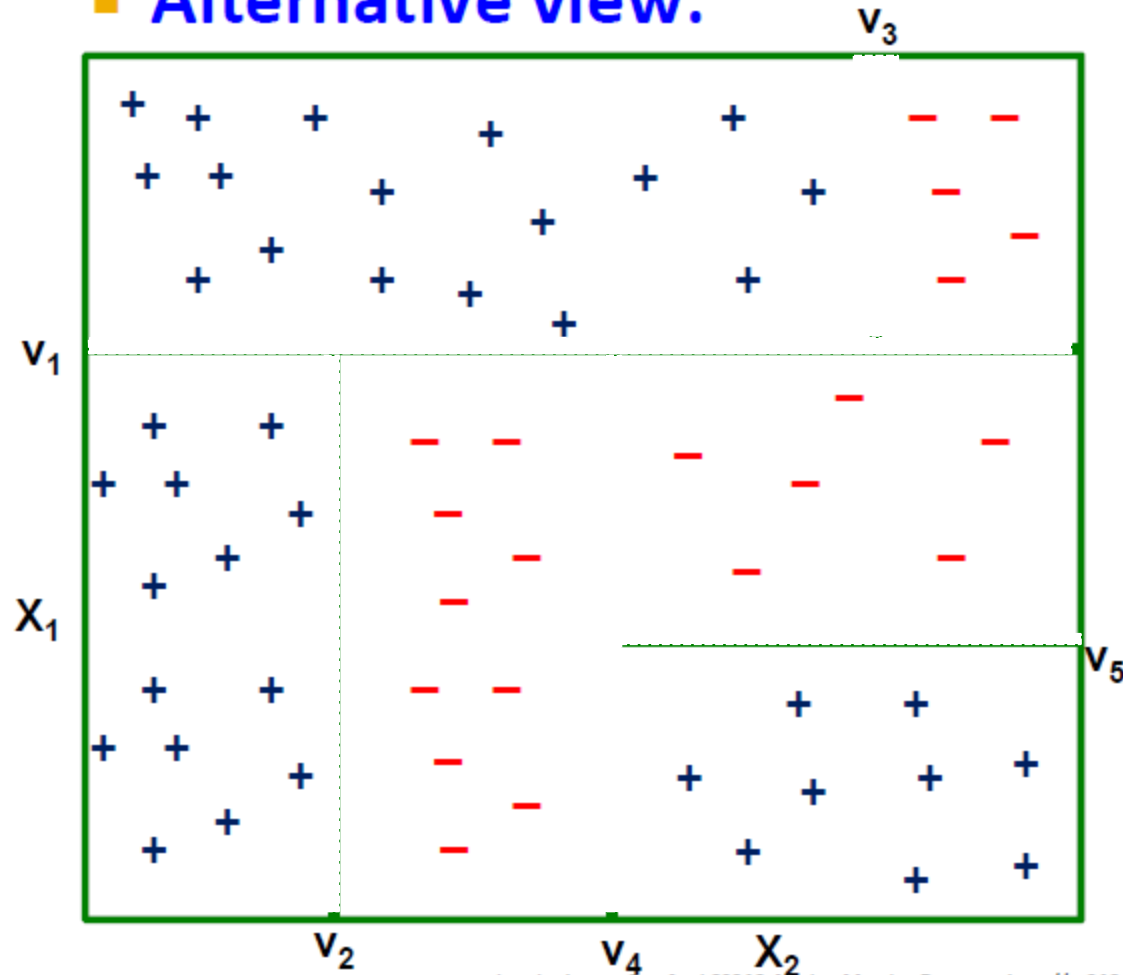
# How to make predictions?

- **Input:** Example  $x_i$
- **Output:** Predicted  $y_i'$
- “Drop”  $x_i$  down the tree until it hits a leaf node
- Predict the value stored in the leaf that  $x_i$  hits



# Alternative View of a Tree

## ■ Alternative view:

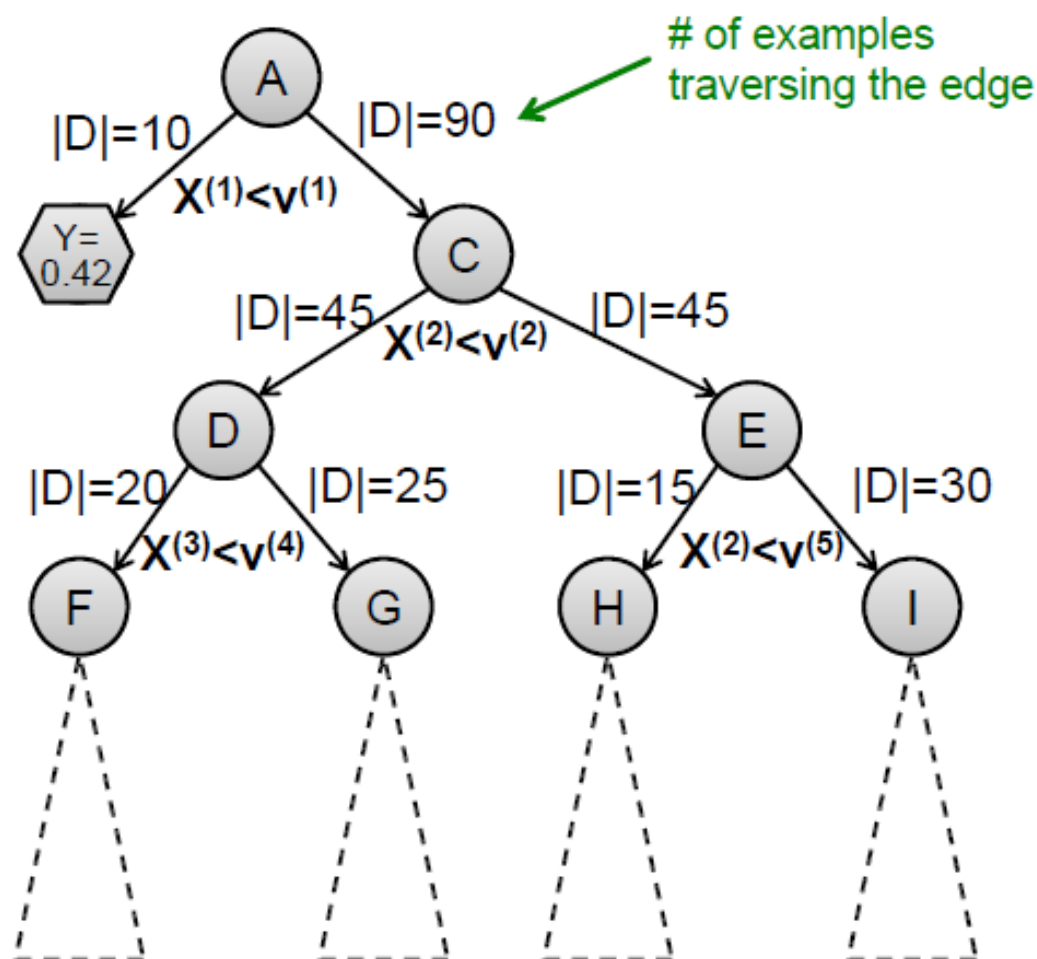


# **HOW TO CONSTRUCT A TREE**



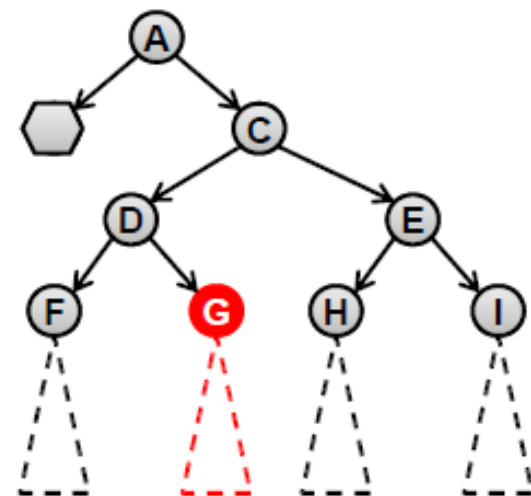
# How to construct a tree?

- Training dataset  $D^*$ ,  $|D^*|=100$  examples



# How to construct a tree?

- Imagine we are currently at some node **G**
  - Let  $D_G$  be the data that reaches **G**
- **There is a decision we have to make: Do we continue building the tree?**
  - **If yes**, which variable and which value do we use for a **split**?
    - Continue building the tree recursively
  - **If not**, how do we make a prediction?
    - We need to build a “**predictor node**”



# How to construct a tree?

---

## Algorithm 1 **BuildSubtree**

---

Require: Node  $n$ , Data  $D \subset D^*$

1:  $(n \rightarrow \text{split}, D_L, D_R) = \text{FindBestSplit}(D)$  (1)

2: if  $\text{StoppingCriteria}(D_L)$  then (2)

3:    $n \rightarrow \text{left\_prediction} = \text{FindPrediction}(D_L)$  (3)

4: else

5:         **BuildSubtree** ( $n \rightarrow \text{left}, D_L$ )

6: if  $\text{StoppingCriteria}(D_R)$  then

7:    $n \rightarrow \text{right\_prediction} = \text{FindPrediction}(D_R)$

8: else

9:         **BuildSubtree** ( $n \rightarrow \text{right}, D_R$ )

---

- Requires at least a single pass over the data!



# Criterion for attribute selection

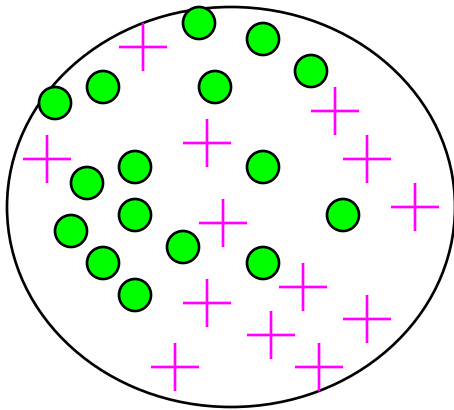
---

- Which is the best split?
  - The one which will result in the smallest tree
  - **Heuristic**: choose the attribute that produces the “purest” nodes
- Need a good measure of purity!
  - Maximal when?
  - Minimal when?

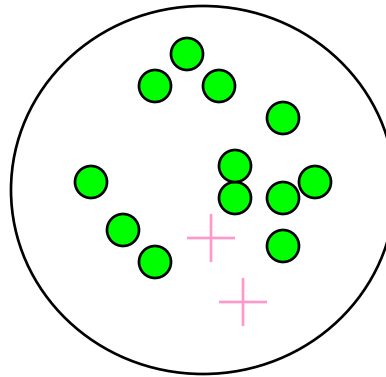


# Impurity

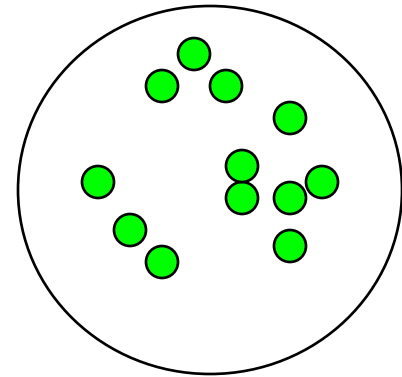
**Very impure group**



**Less impure**



**Minimum impurity**



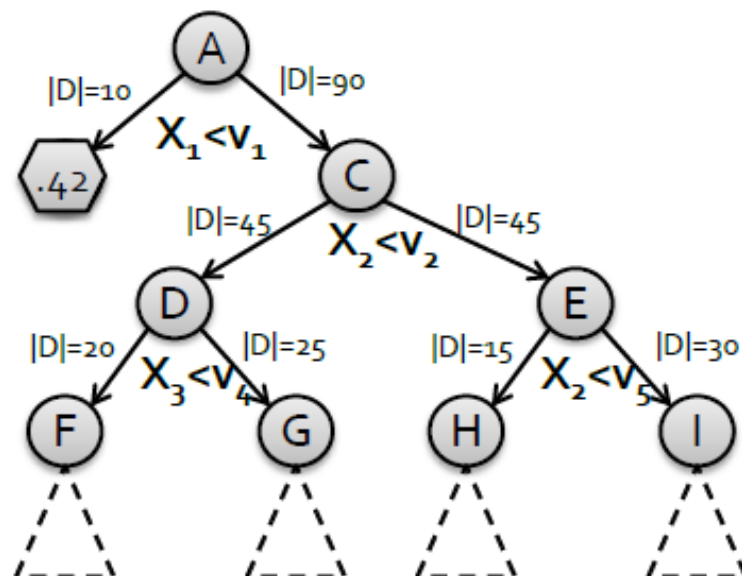
# How to construct a tree?

(1) How to split? Pick attribute & value that optimizes some criterion

- **Classification:**

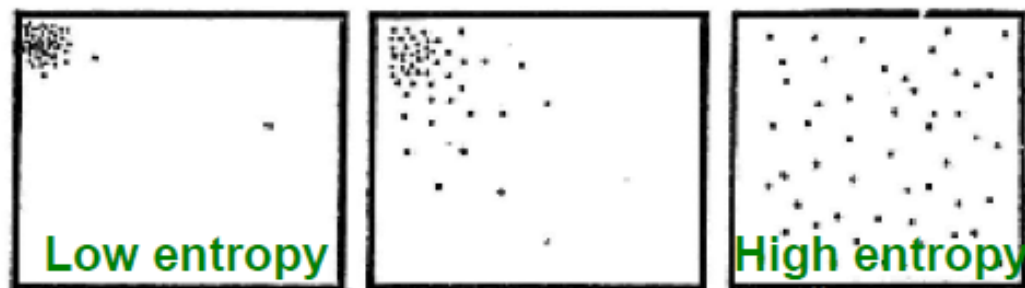
- Information Gain**

- Measures how much a given attribute  $X$  tells us about the class  $Y$
- $IG(Y | X)$  : We must transmit  $Y$  over a binary link. How many bits on average would it save us if both ends of the line knew  $X$ ?



# Why Information Gain? Entropy

- **Entropy:** What's the smallest possible number of bits, on average, per symbol, needed to transmit a stream of symbols drawn from  $\mathbf{X}$ 's distribution?
- **The entropy of  $\mathbf{X}$ :**  $H(\mathbf{X}) = -\sum_{j=1}^m p_j \log p_j$ 
  - “High Entropy”:  $\mathbf{X}$  is from a uniform (boring) distribution
    - A histogram of the frequency distribution of values of  $\mathbf{X}$  is **flat**
  - “Low Entropy”:  $\mathbf{X}$  is from varied (peaks and valleys) distribution
    - A histogram of the frequency distribution of values of  $\mathbf{X}$  would have many lows and one or two highs



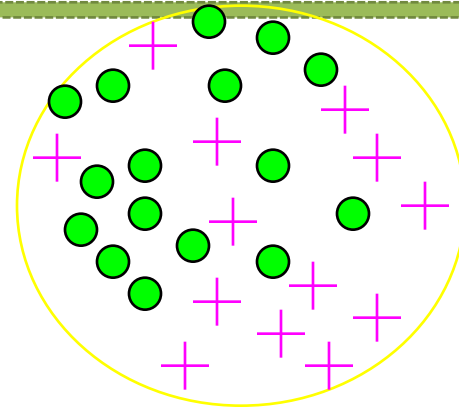


# Entropy: a common way to measure impurity

- Entropy = 
$$\sum_i -p_i \log_2 p_i$$

$p_i$  is the probability of class  $i$

Compute it as the proportion of class  $i$  in the set.



- Entropy comes from information theory. The higher the entropy the more the information content.

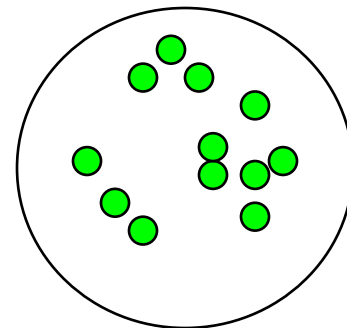
What does that mean for learning from examples?



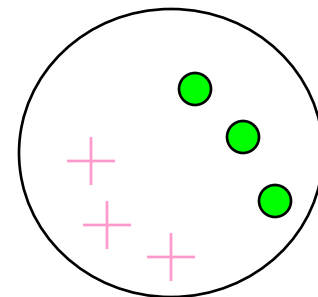
# 2-Class Cases:

- What is the entropy of a group in which all examples belong to the same class?
  - $\text{entropy} = -1 \log_2 1 = 0$
- What is the entropy of a group with 50% in either class?
  - $\text{entropy} = -0.5 \log_2 0.5 - 0.5 \log_2 0.5 = 1$

**Minimum impurity**



**Maximum impurity**





# Exercise

- Compute the entropy of the following distribution
  - $\{ 1, 0, 1, 1, 1, 1, 0, 1, 0, 0, 0, 0, 0, 0, 1, 0 \}$



# Exercise

---

- Compute the entropy of the following distribution
  - { blue, blue, red, orange, red, black, red, black, blue, blue, blue, blue, blue, red }



# Information Gain

---

- We want to determine **which attribute** in a given set of training feature vectors is **most useful** for discriminating between the classes to be learned.
- **Information gain** tells us how important a given attribute of the feature vectors is.
- We will use it to decide the ordering of attributes in the nodes of a decision tree.

# Why Information Gain? Entropy

- Suppose I want to predict  $Y$  and I have input  $X$ 
  - $X$  = College Major
  - $Y$  = Likes “Gladiator”

$X$	$Y$
Math	Yes
History	No
CS	Yes
Math	No
Math	No
CS	Yes
Math	Yes
History	No

- From this data we estimate
  - $P(Y = \text{Yes}) = 0.5$
  - $P(X = \text{Math} \ \& \ Y = \text{No}) = 0.25$
  - $P(X = \text{Math}) = 0.5$
  - $P(Y = \text{Yes} \mid X = \text{History}) = 0$
- Note:
  - $H(Y) = -\frac{1}{2}\log_2(\frac{1}{2}) - \frac{1}{2}\log_2(\frac{1}{2}) = 1$
  - $H(X) = 1.5$

# Why Information Gain? Entropy

- Suppose I want to predict  $Y$  and I have input  $X$ 
  - $X$  = College Major
  - $Y$  = Likes “Gladiator”

$X$	$Y$
Math	Yes
History	No
CS	Yes
Math	No
Math	No
CS	Yes
Math	Yes
History	No

- **Def: Specific Conditional Entropy**

- $H(Y | X=v)$  = The entropy of  $Y$  among only those records in which  $X$  has value  $v$
- **Example:**
  - $H(Y|X = \text{Math}) = 1$
  - $H(Y|X = \text{History}) = 0$
  - $H(Y|X = \text{CS}) = 0$

# Why Information Gain?

- Suppose I want to predict  $Y$  and I have input  $X$ 
  - $X$  = College Major
  - $Y$  = Likes “Gladiator”

$X$	$Y$
Math	Yes
History	No
CS	Yes
Math	No
Math	No
CS	Yes
Math	Yes
History	No

- **Def: Conditional Entropy**

- $H(Y | X)$  = The average specific conditional entropy of  $Y$ 
  - = if you choose a record at random what will be the conditional entropy of  $Y$ , conditioned on that row's value of  $X$
  - = Expected number of bits to transmit  $Y$  if both sides will know the value of  $X$
- $= \sum_j P(X = v_j) H(Y | X = v_j)$

# Why Information Gain?

- Suppose I want to predict  $Y$  and I have input  $X$ 
  - $H(Y | X)$  = The average specific conditional entropy of  $Y$

X	Y
Math	Yes
History	No
CS	Yes
Math	No
Math	No
CS	Yes
Math	Yes
History	No

$$= \sum_j P(X = v_j) H(Y | X = v_j)$$

- **Example:**

$v_j$	$P(X=v_j)$	$H(Y X=v_j)$
Math	0.5	1
History	0.25	0
CS	0.25	0

- **So:**  $H(Y | X) = 0.5 * 1 + 0.25 * 0 + 0.25 * 0 = 0.5$



# Why Information Gain?

- Suppose I want to predict  $Y$  and I have input  $X$

- **Def: Information Gain**

- $IG(Y|X)$  = I must transmit  $Y$ . **How many bits on average would it save me if both ends of the line knew  $X$ ?**

$$IG(Y|X) = H(Y) - H(Y|X)$$

- **Example:**

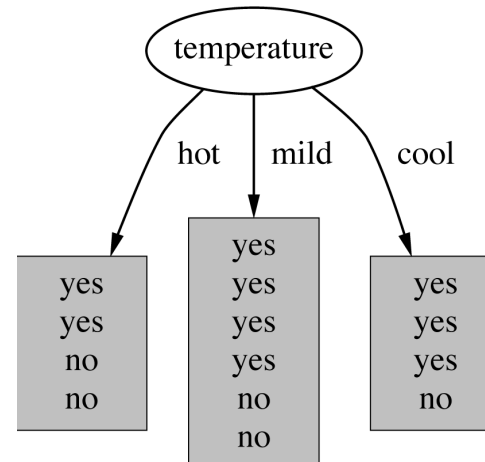
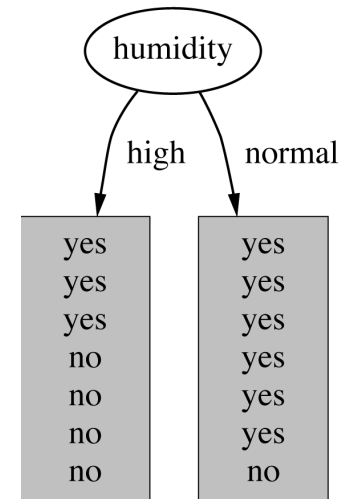
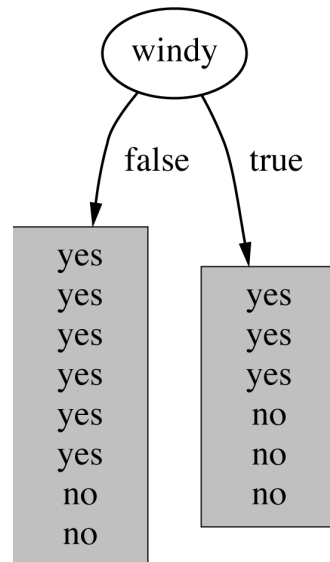
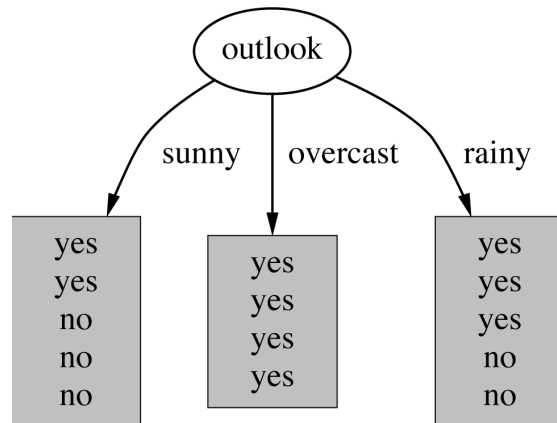
- $H(Y) = 1$
  - $H(Y|X) = 0.5$
  - Thus  $IG(Y|X) = 1 - 0.5 = 0.5$

X	Y
Math	Yes
History	No
CS	Yes
Math	No
Math	No
CS	Yes
Math	Yes
History	No

# What is Information Gain used for?

- Suppose you are trying to predict whether someone is going live past 80 years
- From historical data you might find:
  - $IG(\text{LongLife} \mid \text{HairColor}) = 0.01$
  - $IG(\text{LongLife} \mid \text{Smoker}) = 0.3$
  - $IG(\text{LongLife} \mid \text{Gender}) = 0.25$
  - $IG(\text{LongLife} \mid \text{LastDigitOfSSN}) = 0.00001$
- IG tells us how much information about  $Y$  is contained in  $X$ 
  - So attribute  $X$  that has high  $IG(Y \mid X)$  is a good split!

# Which attribute to select?

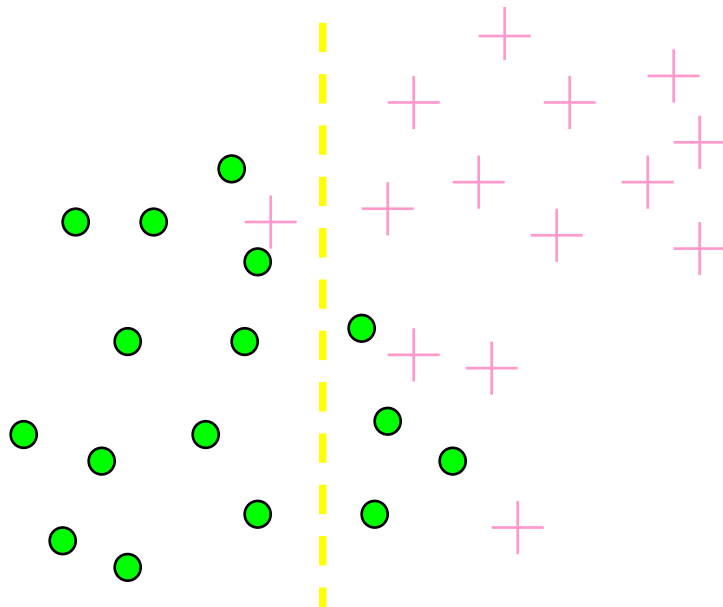




# Information Gain

Which test is more informative?

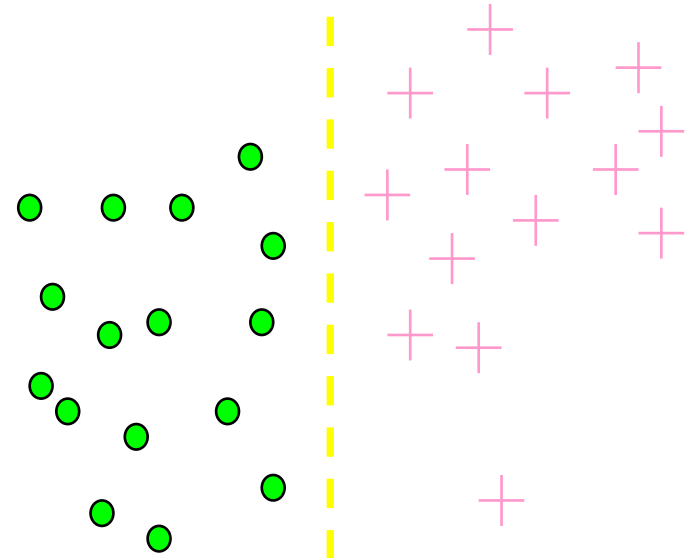
**Split over whether  
Balance exceeds 50K**



Less or equal 50K

Over 50K

**Split over whether  
applicant is employed**



Unemployed

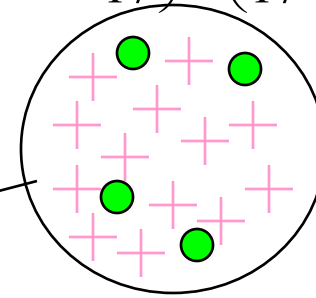
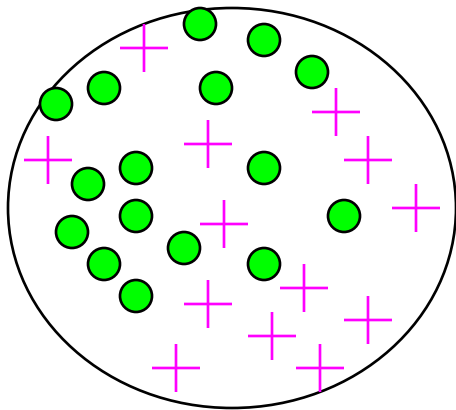
Employed

# Exercise

**Information Gain** = entropy(parent) – [average entropy(children)]

child entropy  $-\left(\frac{13}{17} \cdot \log_2 \frac{13}{17}\right) - \left(\frac{4}{17} \cdot \log_2 \frac{4}{17}\right) = 0.787$

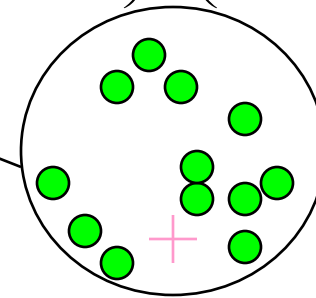
Entire population (30 instances)



17 instances

child entropy  $-\left(\frac{1}{13} \cdot \log_2 \frac{1}{13}\right) - \left(\frac{12}{13} \cdot \log_2 \frac{12}{13}\right) = 0.391$

parent entropy  $-\left(\frac{14}{30} \cdot \log_2 \frac{14}{30}\right) - \left(\frac{16}{30} \cdot \log_2 \frac{16}{30}\right) = 0.996$



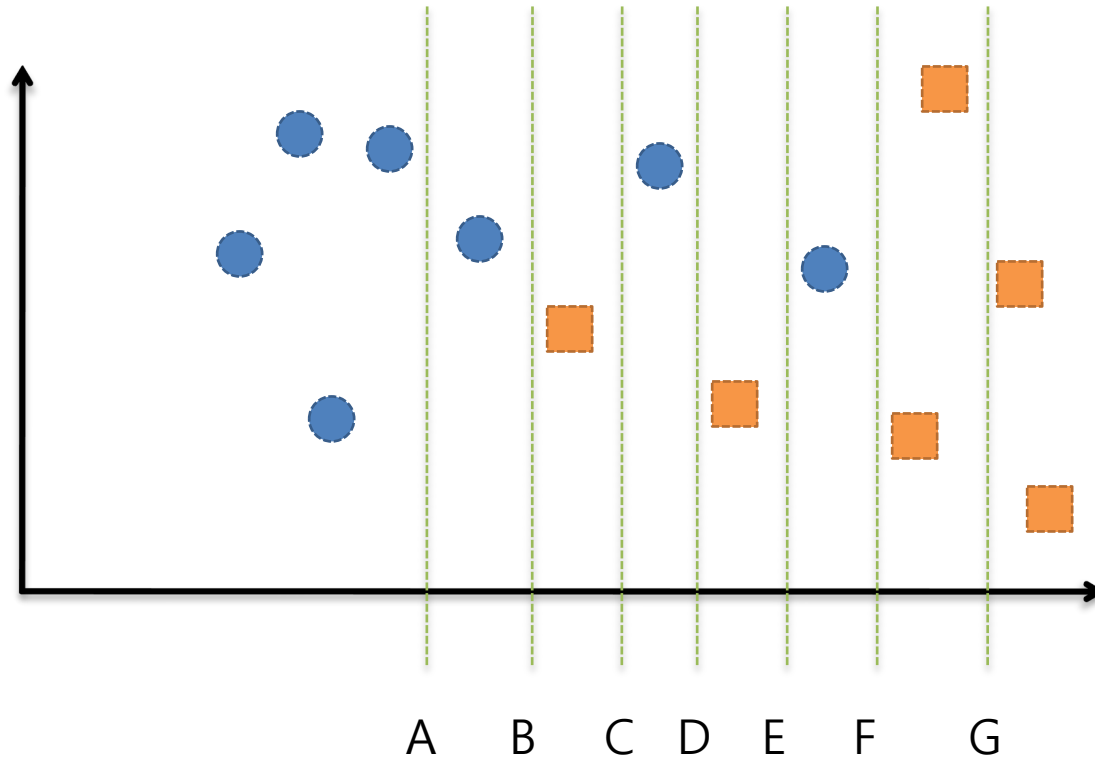
13 instances

(Weighted) Average Entropy of Children =  $\left(\frac{17}{30} \cdot 0.787\right) + \left(\frac{13}{30} \cdot 0.391\right) = 0.615$

**Information Gain = 0.996 - 0.615 = 0.38**



# Exercise



**Which one is the best split point?**

# **DECISION TREE BUILDING ALGORITHM**



# Decision Tree Algorithm

- A decision tree is created in two phases:
  - Building Phase
    - Recursively split nodes using best splitting attribute for node until all the examples in each node belong to one class
  - Pruning Phase
    - Prune leaf nodes recursively to prevent overfitting
    - Smaller imperfect decision tree generally achieves better accuracy





# Building Phase

---

- General tree-growth algorithm (binary tree)
  - Build a tree recursively
- **Partition**(Data S)
- If (all points in S are of the same class) then return;
- for each attribute A do
- evaluate splits on attribute A;
- Use **best split** to partition S into S1 and S2;
- *Partition*(S1);
- *Partition*(S2);



# Summary

## How to construct a tree?

---

### Algorithm 1 BuildSubtree

---

Require: Node  $n$ , Data  $D \subset D^*$

1:  $(n \rightarrow \text{split}, D_L, D_R) = \text{FindBestSplit}(D)$  (1)

2: if  $\text{StoppingCriteria}(D_L)$  then (2)

3:    $n \rightarrow \text{left\_prediction} = \text{FindPrediction}(D_L)$  (3)

4: else

5:       **BuildSubtree** ( $n \rightarrow \text{left}, D_L$ )

6: if  $\text{StoppingCriteria}(D_R)$  then

7:    $n \rightarrow \text{right\_prediction} = \text{FindPrediction}(D_R)$

8: else

9:       **BuildSubtree** ( $n \rightarrow \text{right}, D_R$ )

---

- Requires at least a single pass over the data!