

# Large Scale Machine Learning: Decision Trees

Mining of Massive Datasets  
Leskovec, Rajaraman, and Ullman  
Stanford University

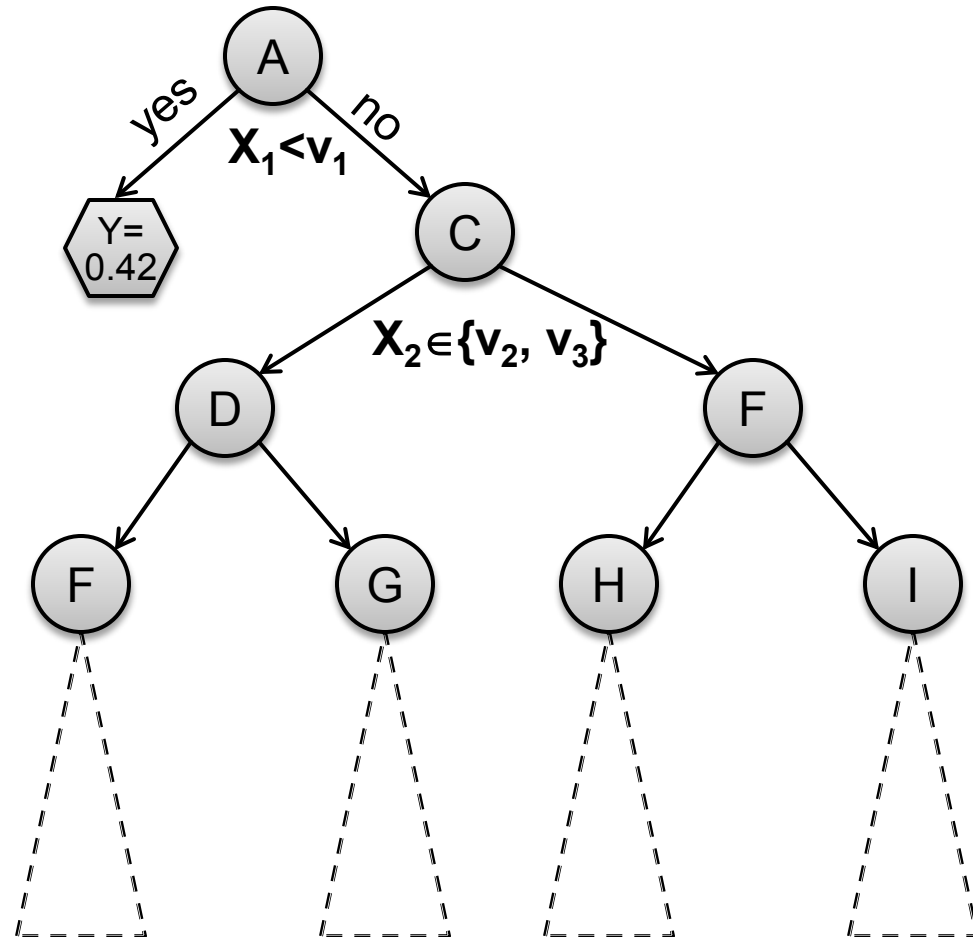


# Decision Tree Learning

- Give one attribute (e.g., wealth), try to predict the value of new people's wealths by means of some of the other available attribute
- **Input attributes:**
  - $d$  features/attributes:  $X_1, X_2, \dots, X_d$
  - Each  $X_j$  has **domain**  $O_j$ 
    - **Categorical:**  $O_j = \{\text{red, blue}\}$
    - **Numerical:**  $H_j = (0, 10)$
  - $Y$  is output variable with domain  $O_Y$ :
    - **Categorical:** Classification, **Numerical:** Regression
- **Data  $D$ :**
  - $n$  examples  $(x_i, y_i)$  where  $x_i$  is a  $d$ -dim feature vector,  $y_i \in O_Y$  is output variable
- **Task:**
  - Given an input data vector  $x$  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



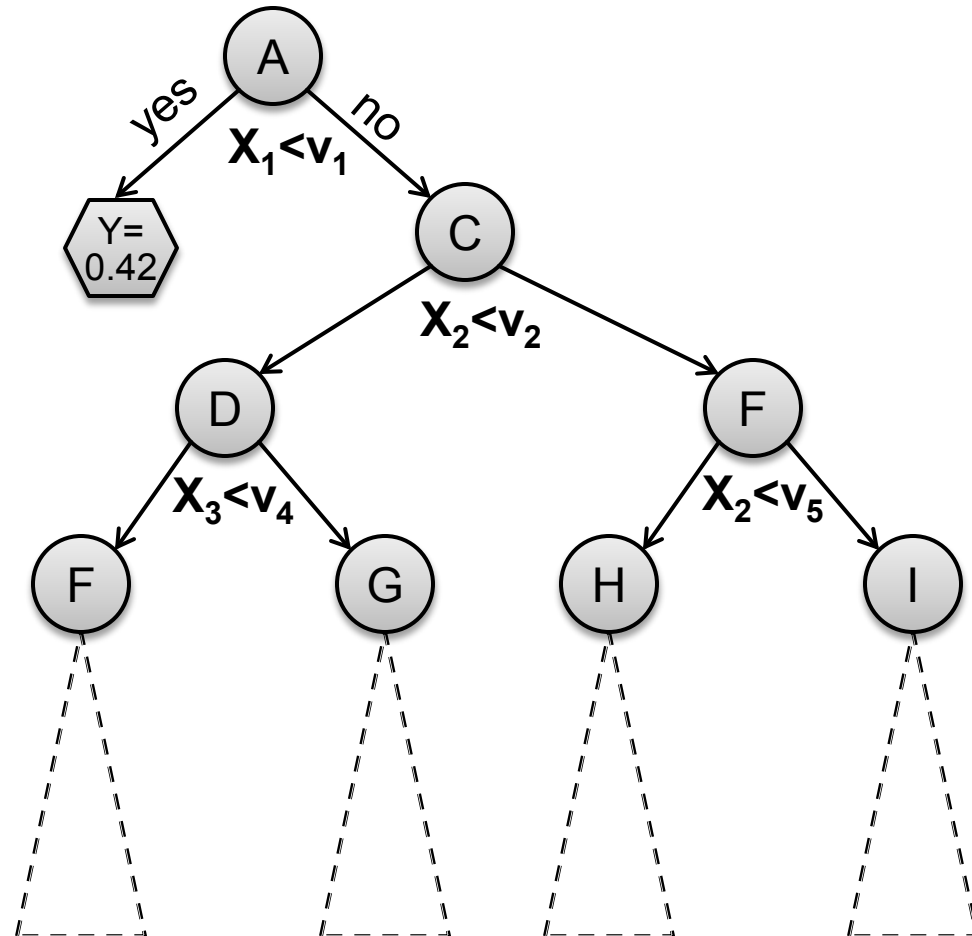
# Decision Trees (1)

## ■ Decision trees:

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

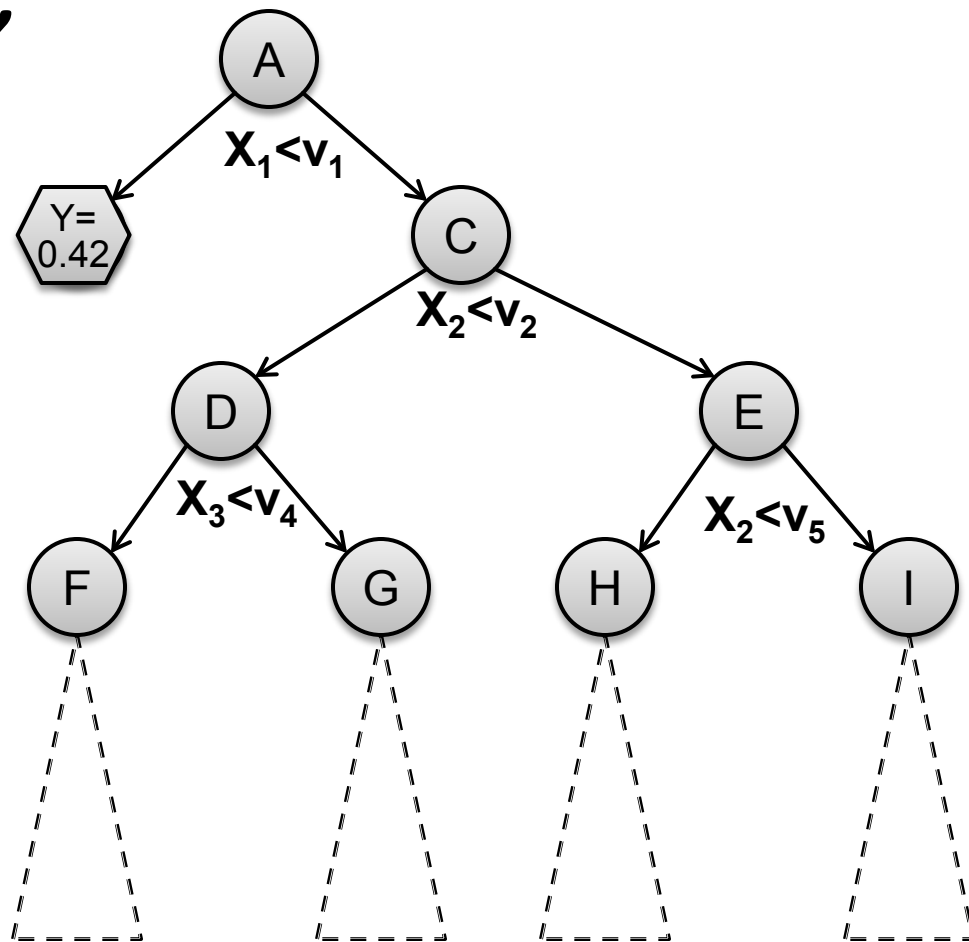
## ■ Lecture today:

- 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



# Decision Trees Vs. SVM

## ■ Alternative view:

