

# CS5228 LECTURE 6: CLASSIFICATION II: TREES & ENSEMBLES

Bryan Hooi

School of Computing

National University of Singapore

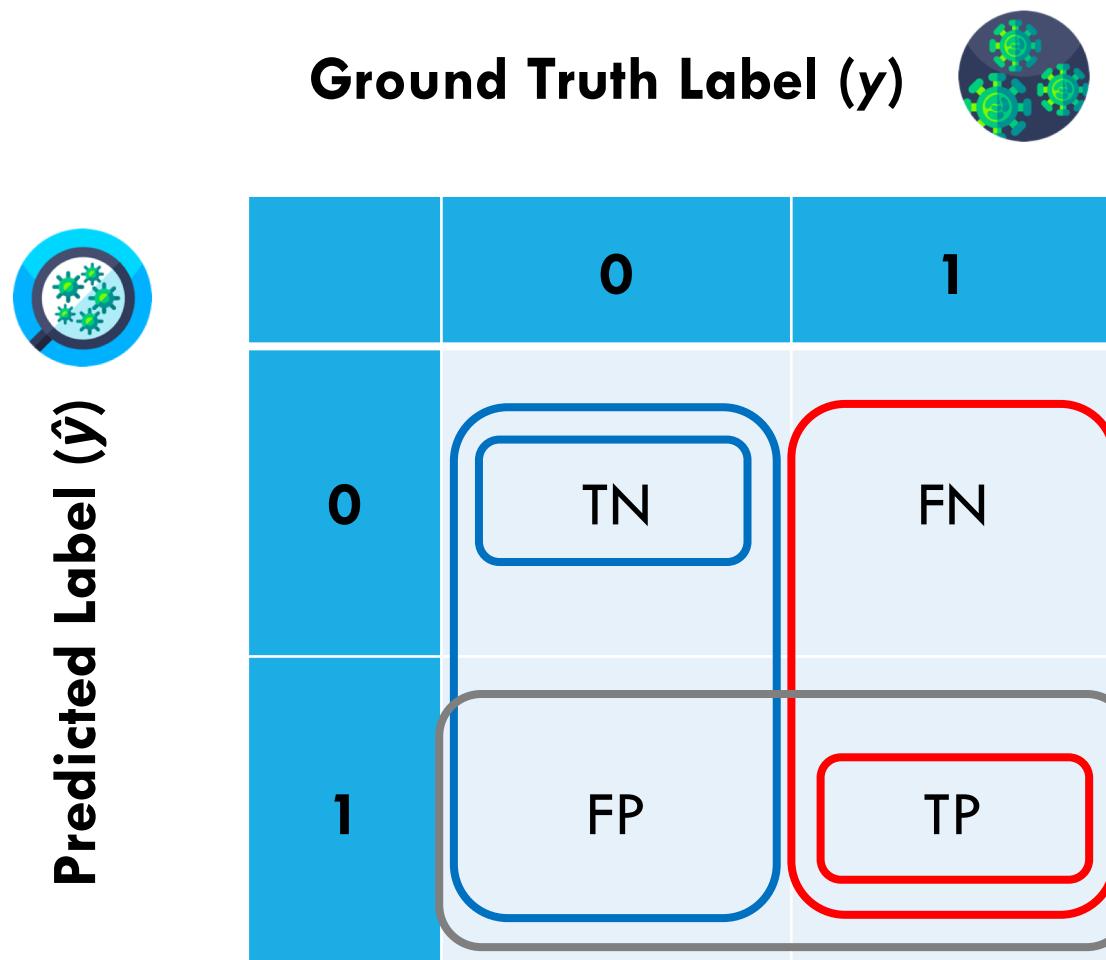
Slide Credit: Wang Wei, Ng See Kiong,  
Wynne Hsu, Chris von der Weth

# ANNOUNCEMENTS

- Assignment 1 is due this Sunday 11.59pm (unless you use late days). You have 8 late days; see the recent announcement for details.
- Assignment 2 will be released this weekend.
- Please finalize your project group in Canvas if you have not done so (deadline is next Friday; email course staff for any difficulties)

Week	Date	Topics	Tutorials	Important Dates
1	13 Jan	Introduction		
2	20 Jan	No class (public holiday)		
3	27 Jan	Clustering I	Tutorial 1	
4	3 Feb	Clustering II		Release A1 + project
5	10 Feb	Association Rules	Tutorial 2	
6	17 Feb	Regression & Classification I		
Recess		No class		
7	3 Mar	Regression & Classification II	Tutorial 3	A1 due (Sunday 11.59pm), release A2
8	10 Mar	Regression & Classification III		
9	17 Mar	Recommender Systems	Tutorial 4	
10	24 Mar	Graph Mining		
11	31 Mar	Data Stream Mining	Tutorial 5	A2 due (Sunday 11.59pm)
12	7 Apr	No class (public holiday)		
13	14 Apr	Review & Outlook		Project due (Sunday 11.59pm)

# BINARY CLASSIFICATION METRICS



**Sensitivity:** same as recall (fraction of actual positives that are correctly identified)

$$= \text{TP} / (\text{TP} + \text{FN})$$

**Specificity:** fraction of actual negatives that are correctly identified

$$= \text{TN} / (\text{TN} + \text{FP})$$

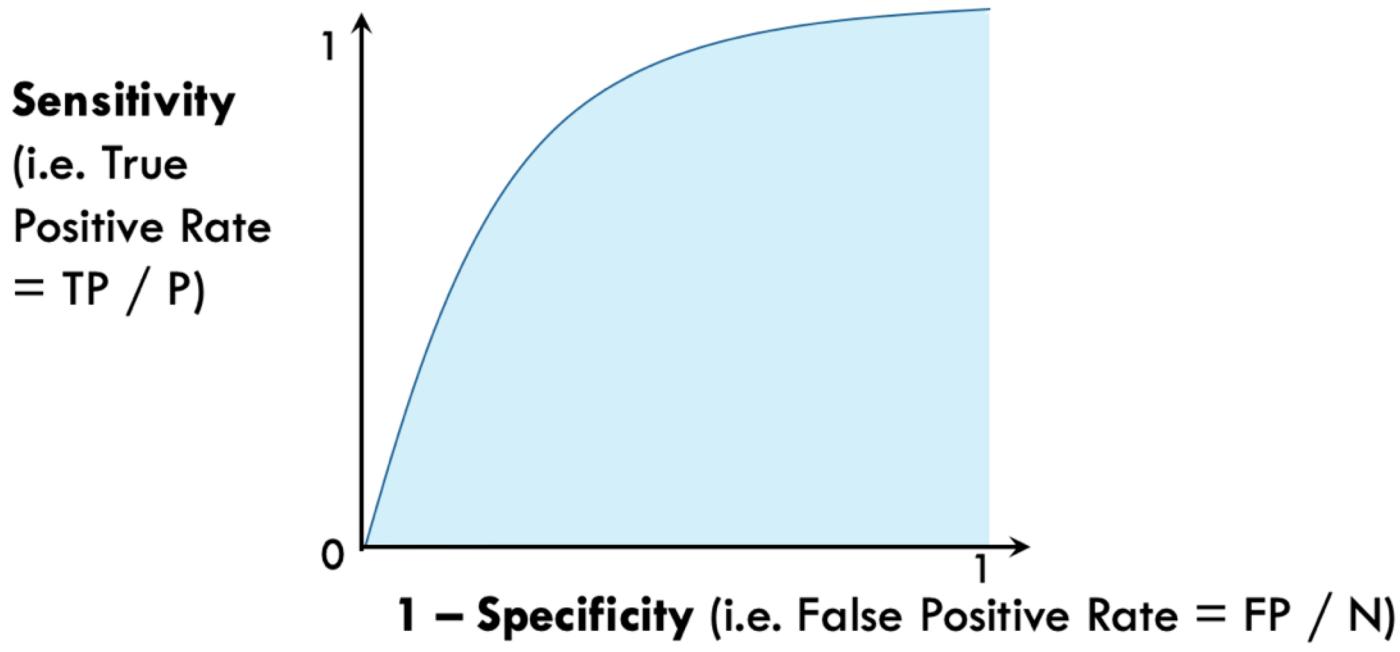
**Interpretation:**

**Sensitivity** = accuracy among positives ( $y=1$ )

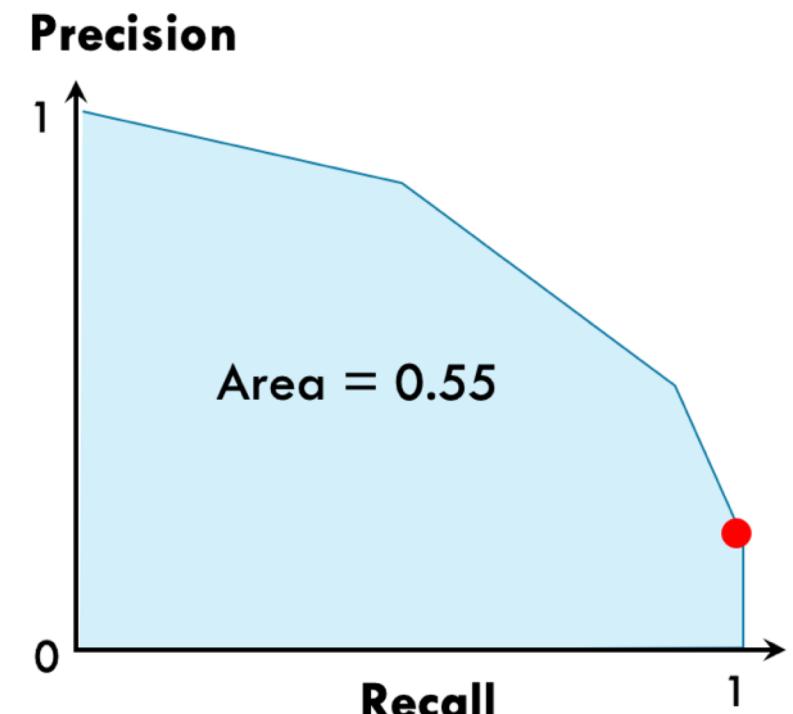
**Specificity** = accuracy among negatives ( $y=0$ )

Precision = accuracy among detected ( $\hat{y}=1$ )

# REVIEW: AUC, AUPRC



**Area Under ROC Curve (AUC)**



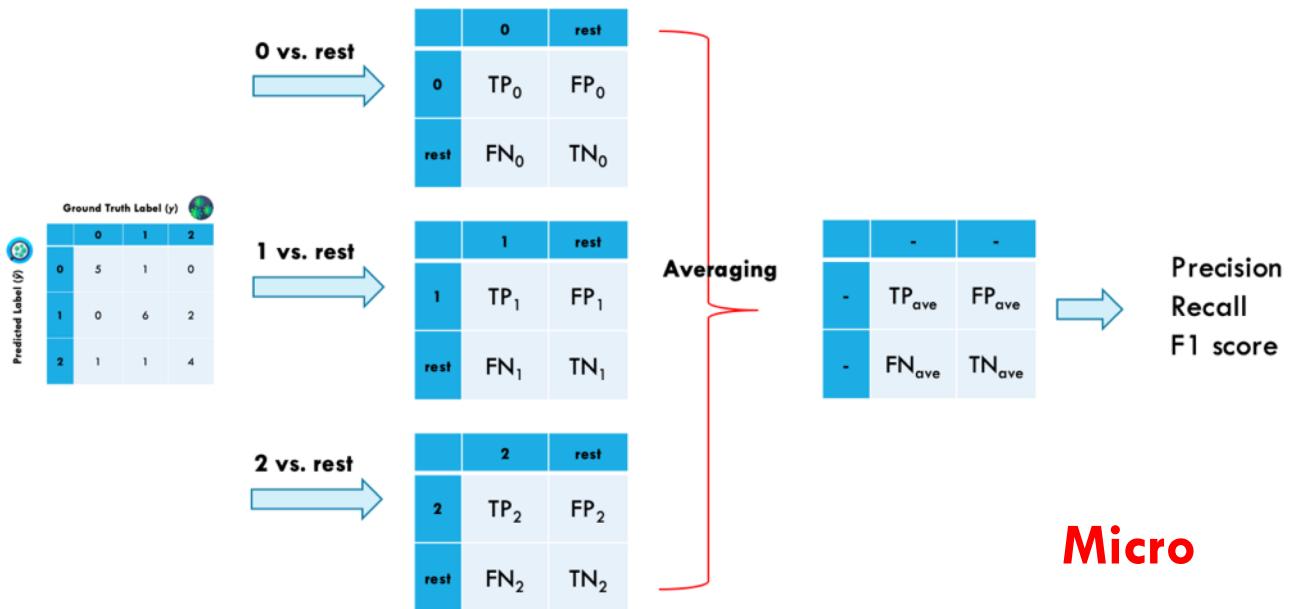
**Area Under Precision Recall Curve (AUPRC)**

# REVIEW: MICRO VS MACRO PRECISION / RECALL / F1

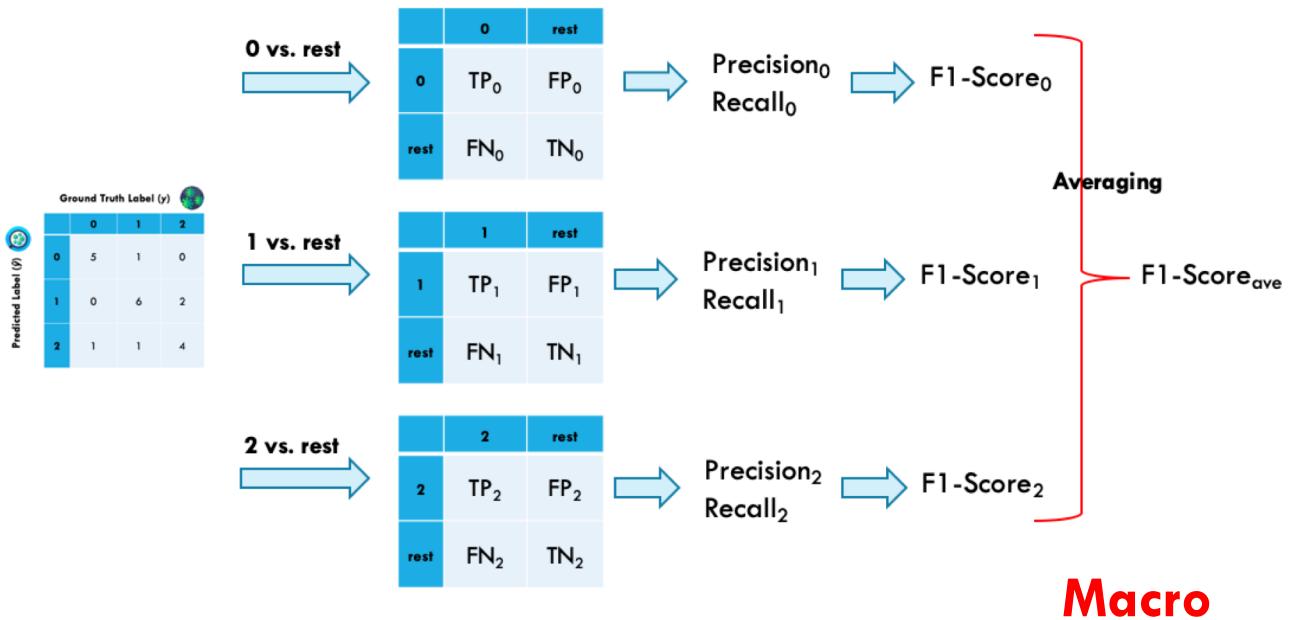
Both based on **one vs rest** confusion matrices

**Micro:** averaging takes place over (TP, TN, FP, FN)

**Macro:** averaging takes place over the statistic itself (e.g. F1-score)



**Micro**

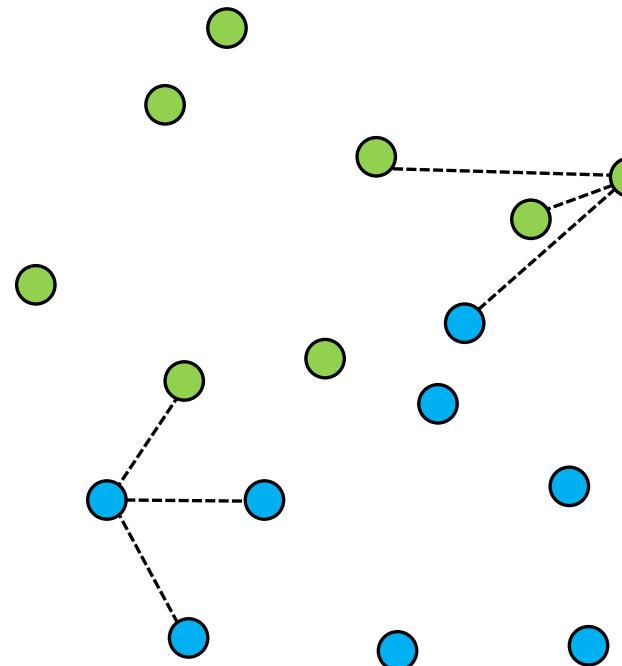


**Macro**

# REVIEW: K NEAREST NEIGHBORS ALGORITHM

For each test point  $x$ :

- Compute its distances to all training instances:  $d_1, \dots, d_n$
- **Majority Vote:** Output most frequent class among the  $k$  nearest neighbors of  $x$



$$k = 3$$

# CLASSIFICATION OVERVIEW

1. Problem Setup
2. Evaluating Classifiers
3. Nearest Neighbor Methods

4. **Trees and Ensembles**
  - a) Decision Trees
  - b) Bagging and Random Forests
  - c) Boosting and Gradient Boosting Machines

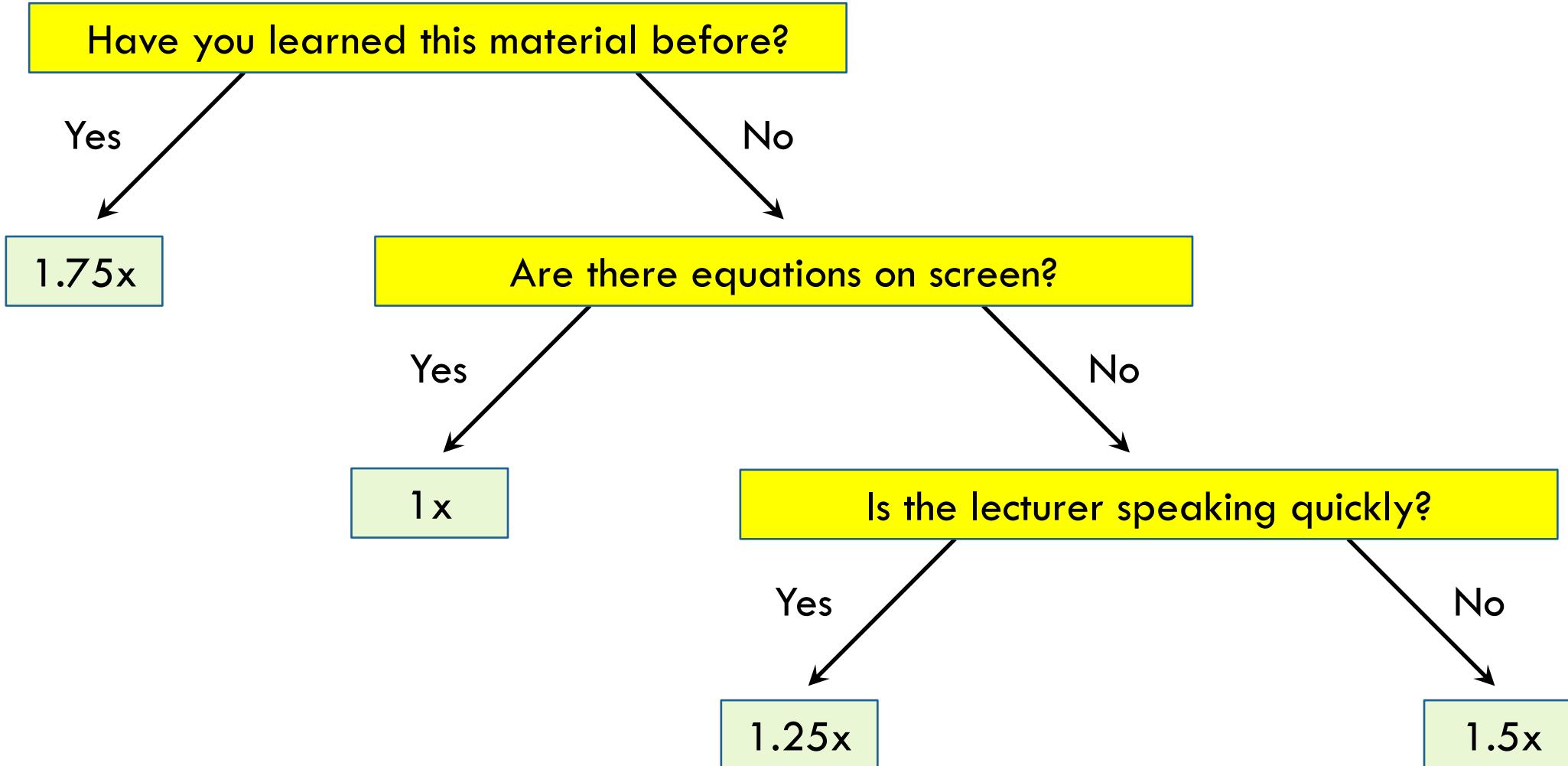
Today's  
Lecture

5. Logistic Regression
6. Deep Learning

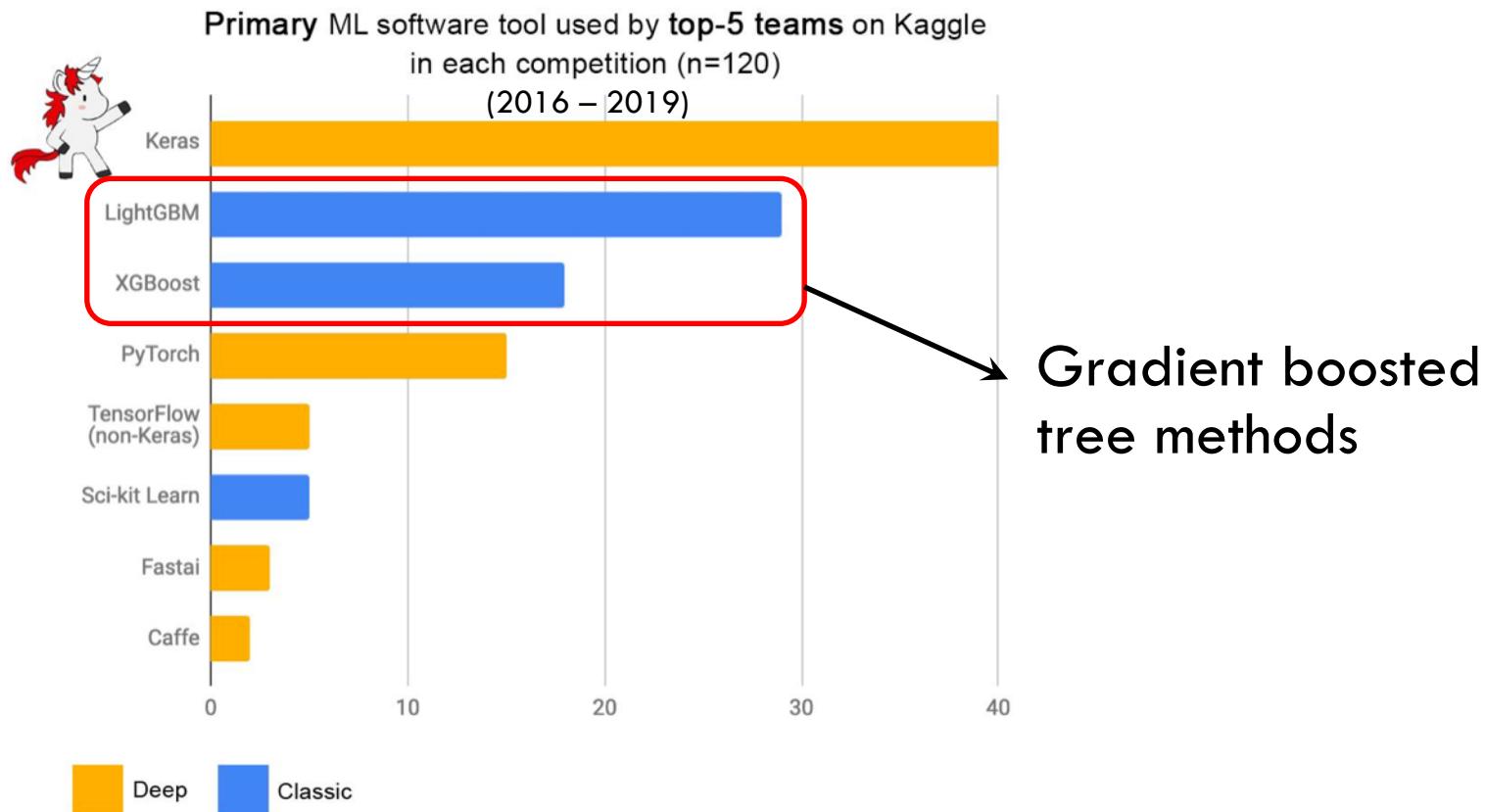
# CLASSIFICATION OVERVIEW

1. Problem Setup
2. Evaluating Classifiers
3. Nearest Neighbor Methods
4. **Trees and Ensembles**
  - a) Decision Trees
  - b) Bagging and Random Forests
  - c) Boosting and Gradient Boosting Machines
5. Logistic Regression
6. Deep Learning

# AT WHAT SPEED SHOULD I WATCH THE LECTURE VIDEO?



# TREE-BASED METHODS ACHIEVE STATE OF THE ART RESULTS

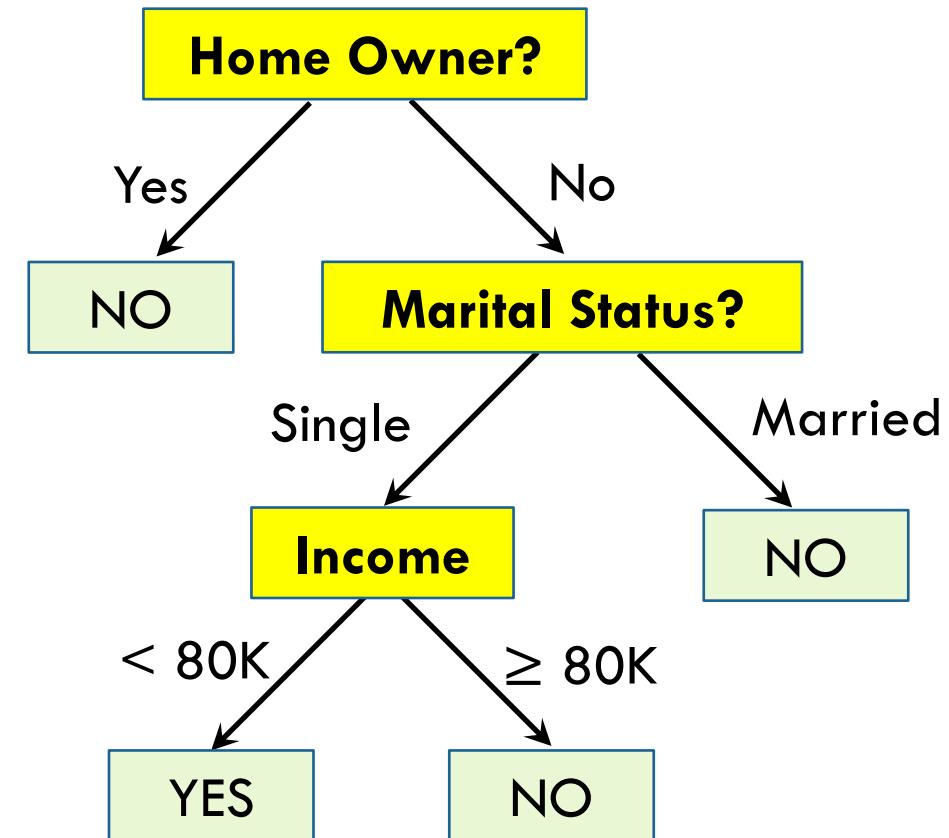
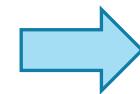


Source: Francois Chollet

Ke G, et al.. **Lightgbm**: A highly efficient gradient boosting decision tree. NeurIPS 2017  
Chen, Tianqi, and Carlos Guestrin. "**Xgboost**: A scalable tree boosting system." KDD 2016

# WHAT IS A DECISION TREE?

ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower	response
1	Yes	Single	125K	No	categorical
2	No	Married	100K	No	categorical
3	No	Single	170K	No	continuous
4	Yes	Married	120K	No	continuous
5	No	Single	75K	Yes	continuous
6	No	Married	160K	No	continuous
7	No	Single	50K	Yes	continuous

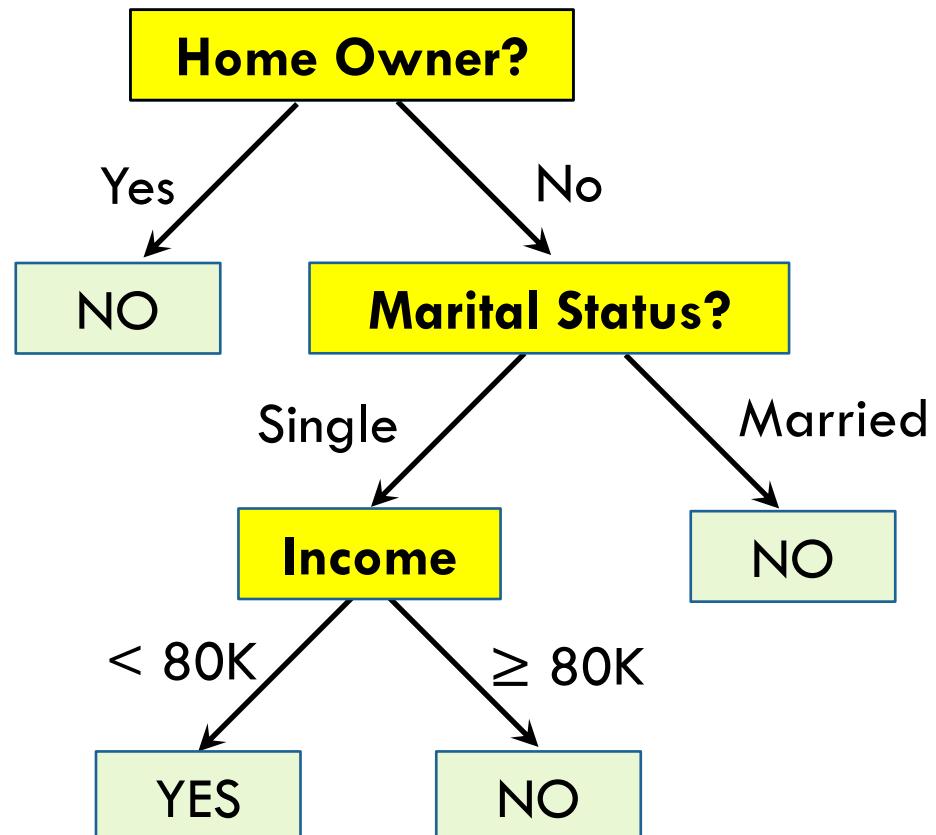


# APPLYING MODEL TO TEST DATA

ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower?
1	No	Married	205K	?

**Test Record**

categorical      categorical      continuous      response

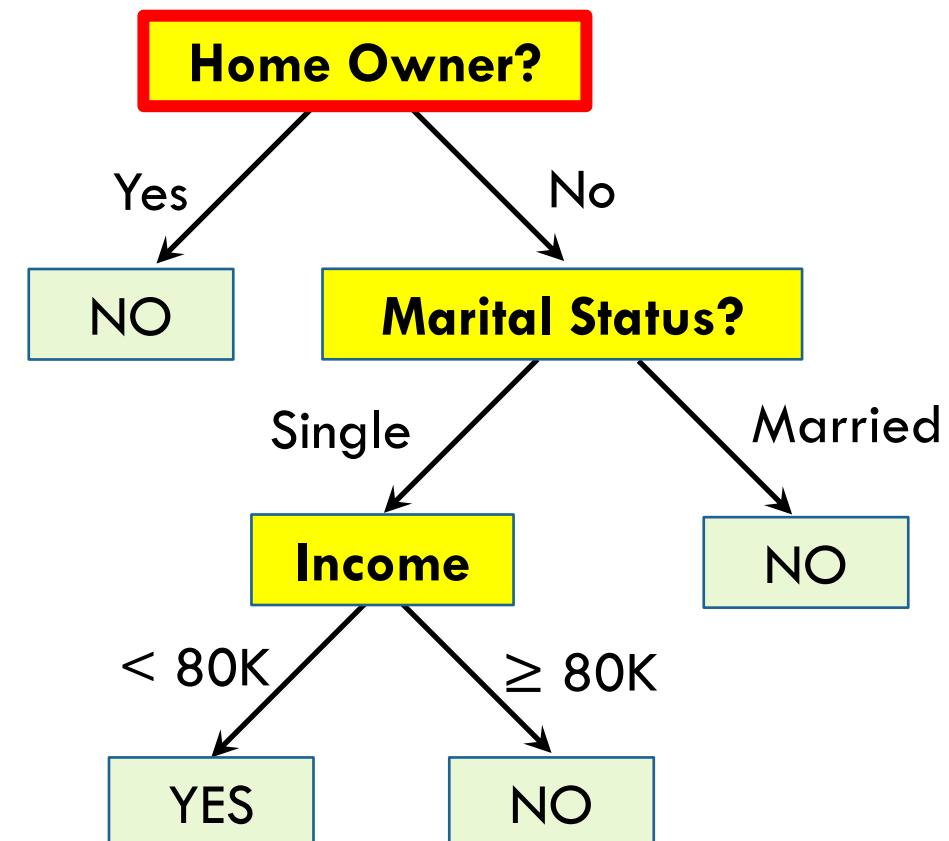


# APPLYING MODEL TO TEST DATA

ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower?
1	No	Married	205K	?

**Test Record**

categorical      categorical      continuous      response

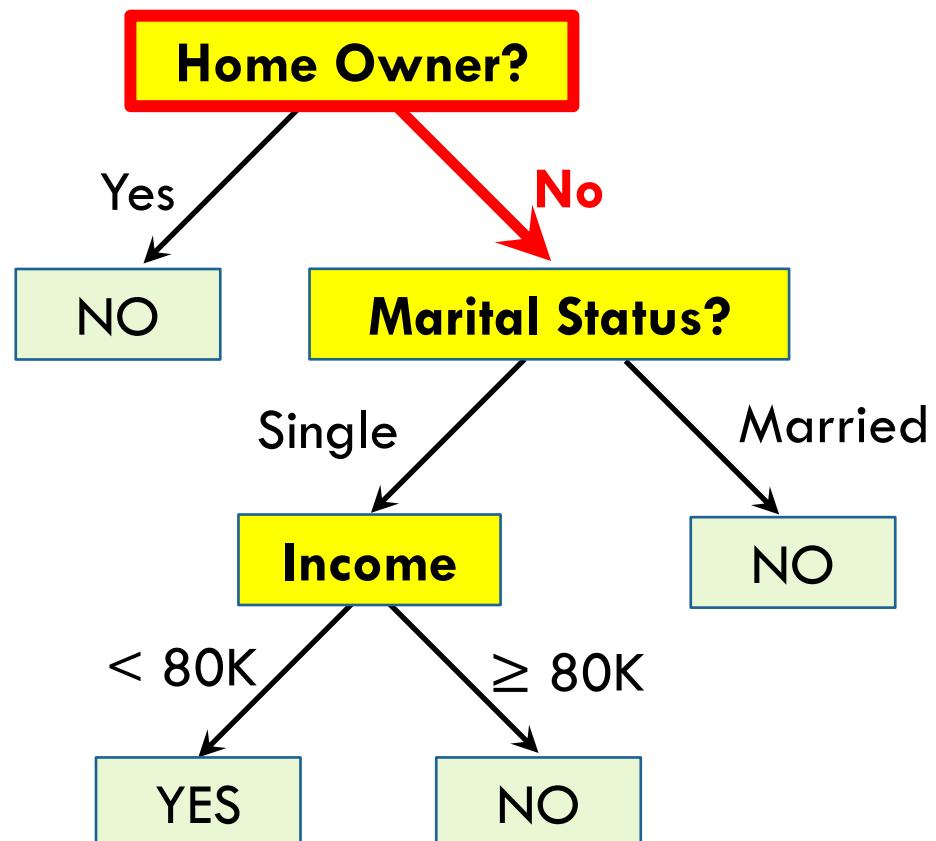


# APPLYING MODEL TO TEST DATA

ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower?
1	No	Married	205K	?

**Test Record**

categorical      categorical      continuous      response

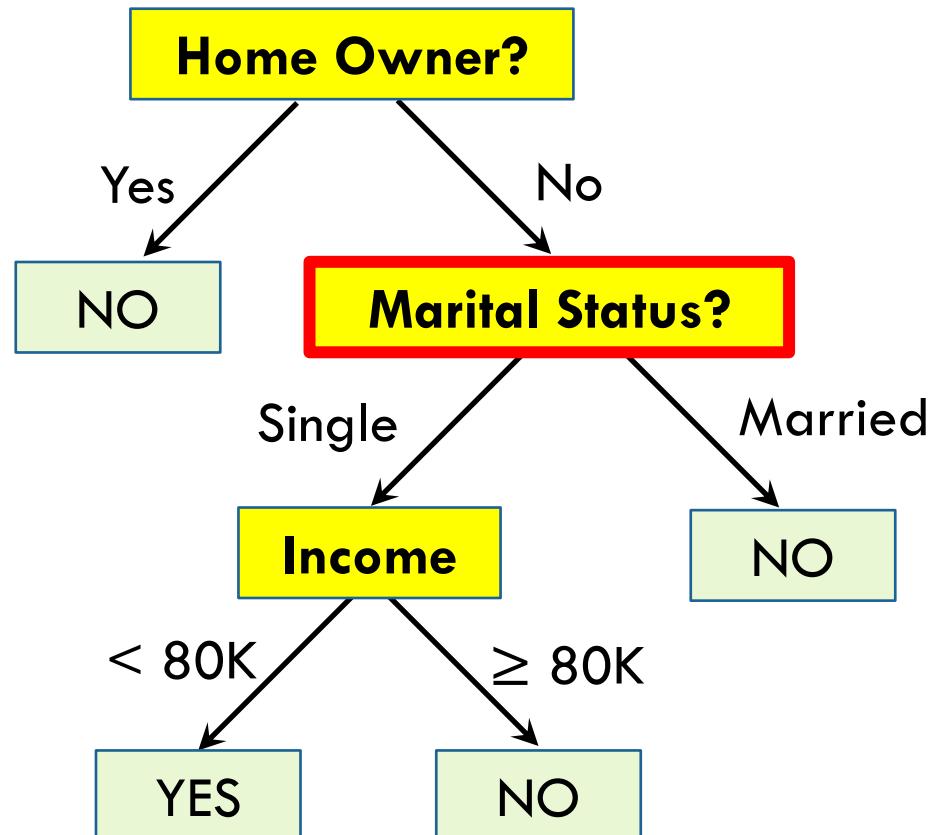


# APPLYING MODEL TO TEST DATA

ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower?
1	No	Married	205K	?

**Test Record**

categorical      categorical      continuous      response

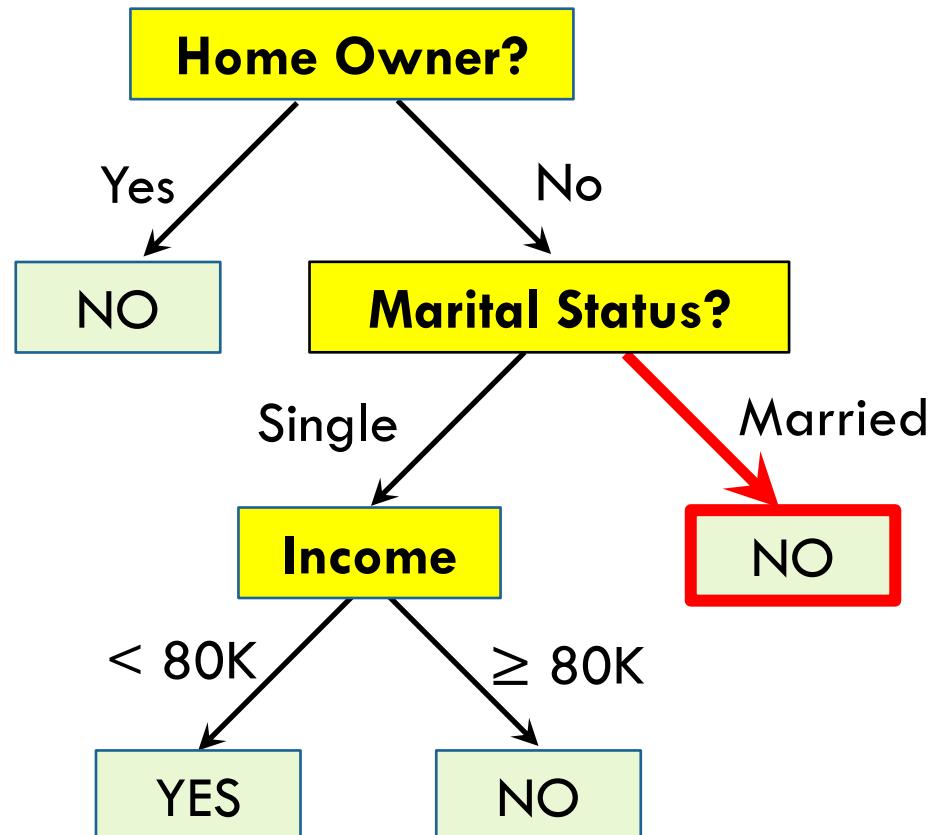


# APPLYING MODEL TO TEST DATA

ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower?
1	No	Married	205K	?

**Test Record**

categorical      categorical      continuous      response

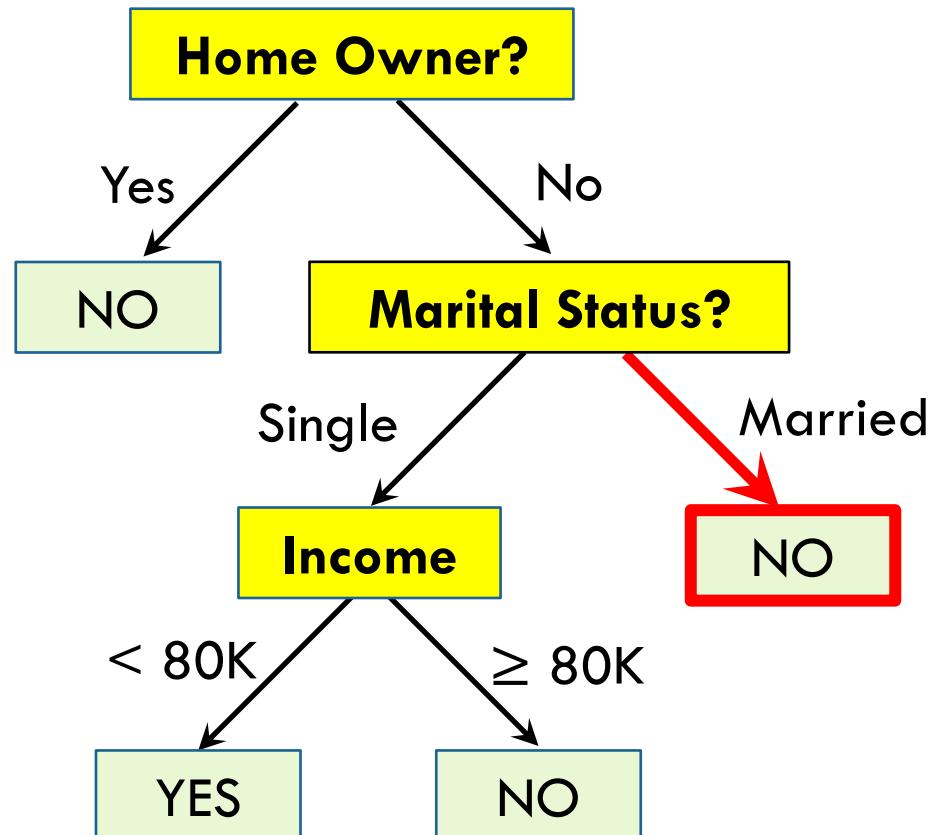


# APPLYING MODEL TO TEST DATA

ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower?
1	No	Married	205K	No

**Test Record**

categorical      categorical      continuous      response



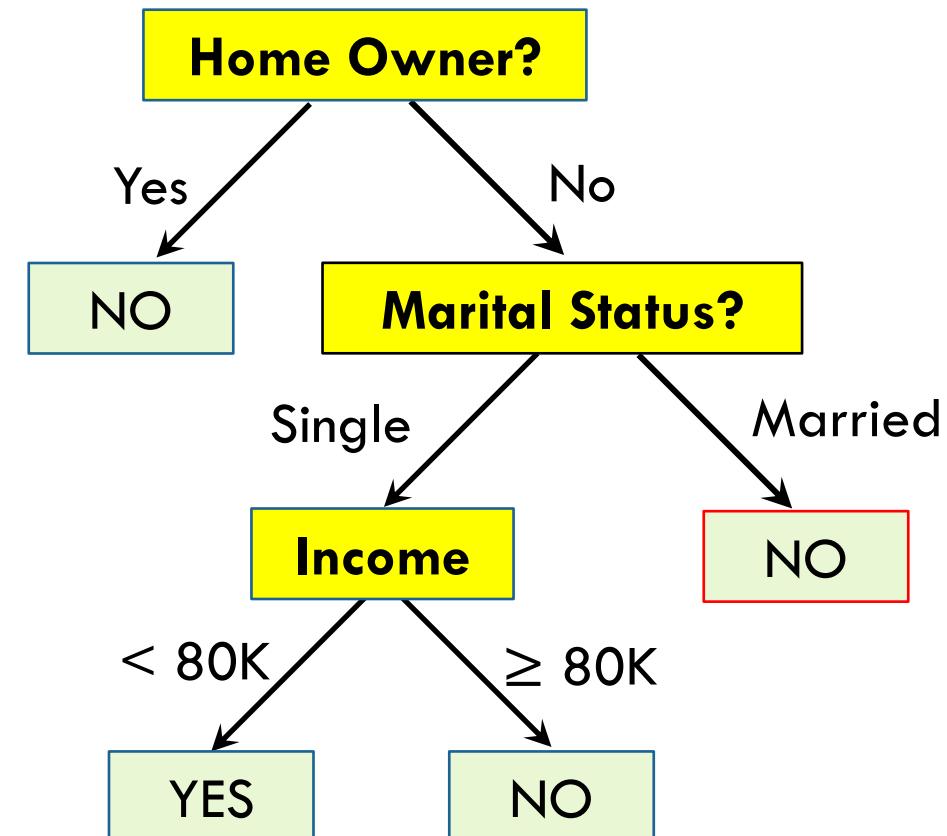
# APPLYING MODEL TO TEST DATA



ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower?
1	No	Single	80K	?

categorical      categorical      continuous      response

**Q:** What is the output of the decision tree on this new test record?



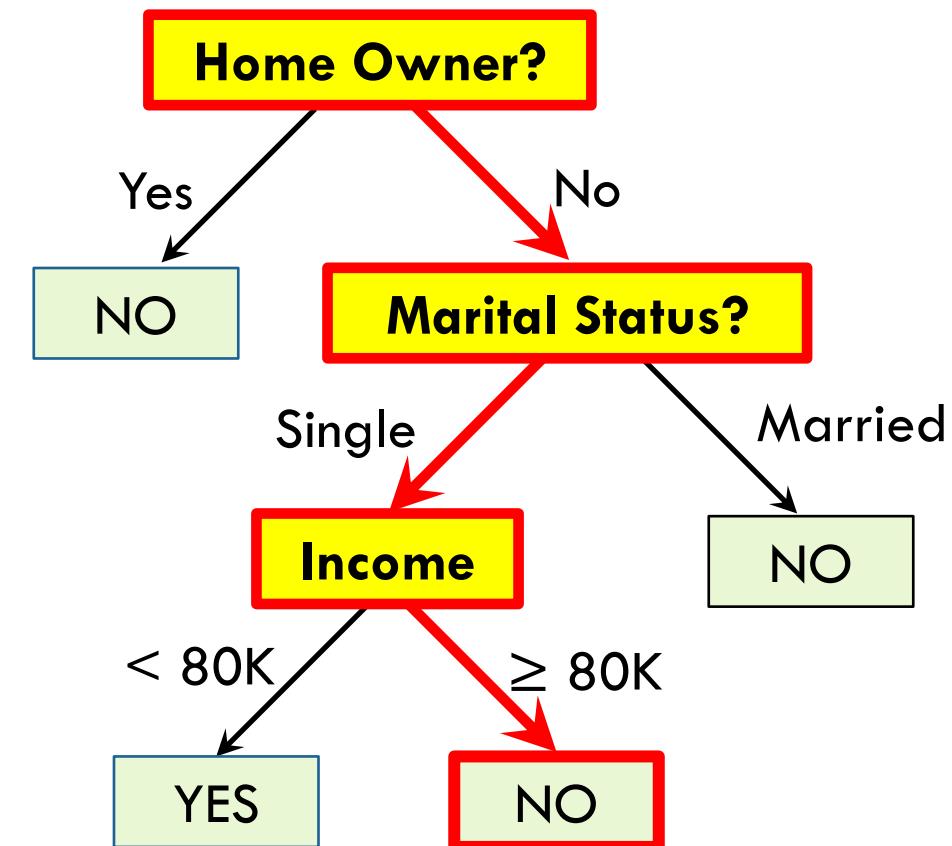
# APPLYING MODEL TO TEST DATA



ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower?
1	No	Single	80K	?

categorical      categorical      continuous      response

A: NO

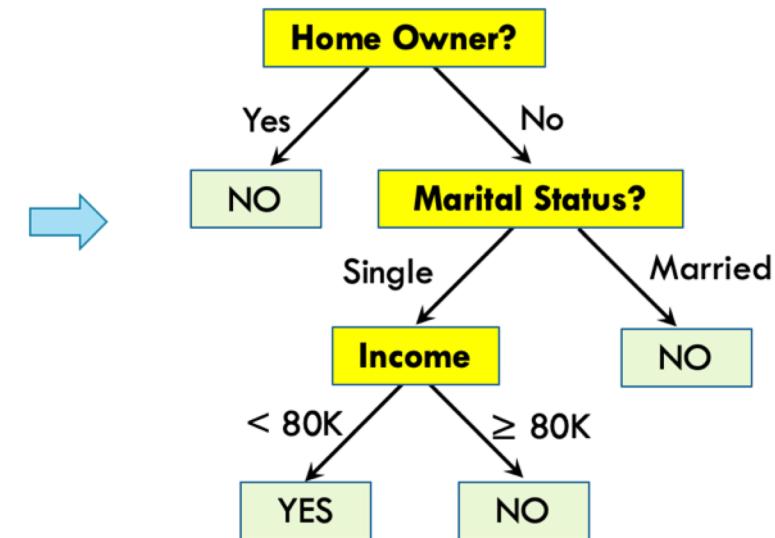


# HOW TO TRAIN DECISION TREE?

Finding the optimal decision tree (for certain notions of optimality) is NP-complete!

Instead, most algorithms use heuristic (i.e. sub-optimal) approaches

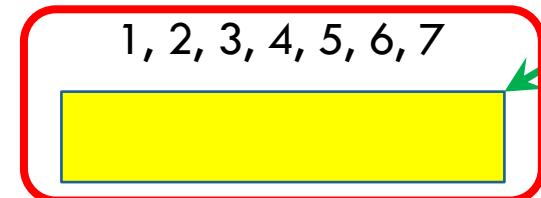
ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	170K	No
4	Yes	Married	120K	No
5	No	Single	75K	Yes
6	No	Married	160K	No
7	No	Single	50K	Yes



# TRAINING DECISION TREE

ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	170K	No
4	Yes	Married	120K	No
5	No	Single	75K	Yes
6	No	Married	160K	No
7	No	Single	50K	Yes

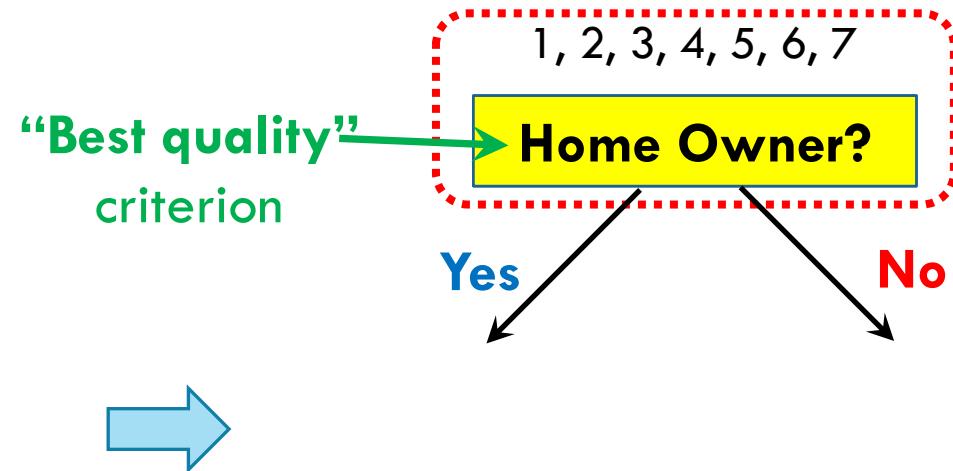
IDs of training data points



1. Initially, the **root node** contains all the training points

# TRAINING DECISION TREE

ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	170K	No
4	Yes	Married	120K	No
5	No	Single	75K	Yes
6	No	Married	160K	No
7	No	Single	50K	Yes

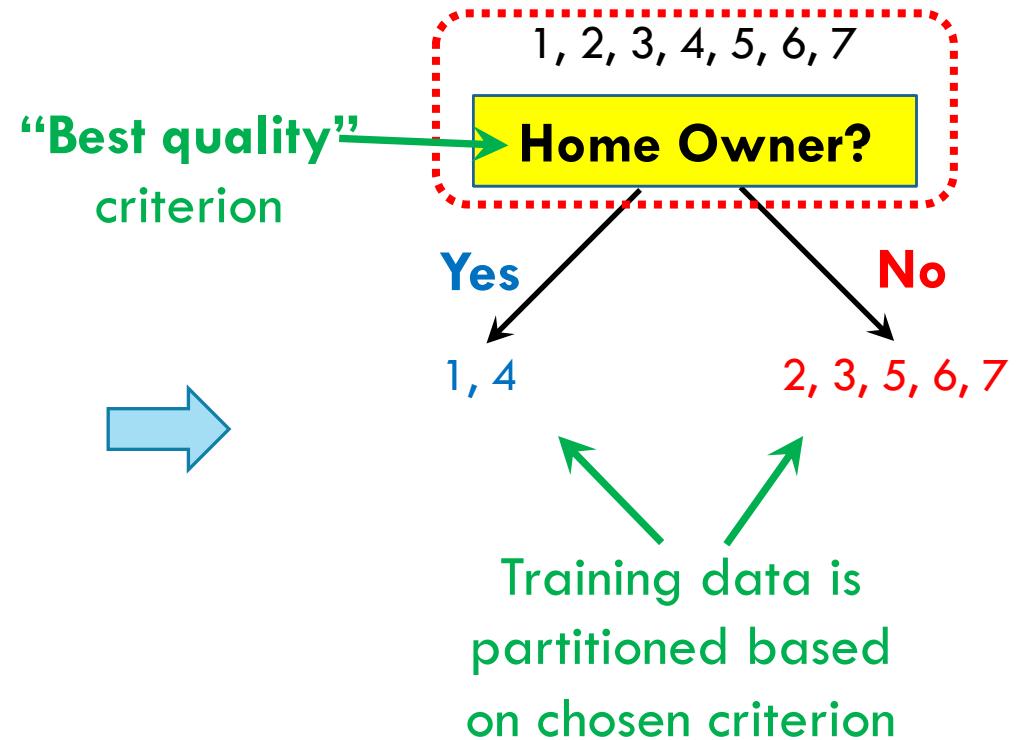


To be defined later

2. Find the “**best quality**” criterion, which partitions the data into 2 subsets

# TRAINING DECISION TREE

ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	170K	No
4	Yes	Married	120K	No
5	No	Single	75K	Yes
6	No	Married	160K	No
7	No	Single	50K	Yes

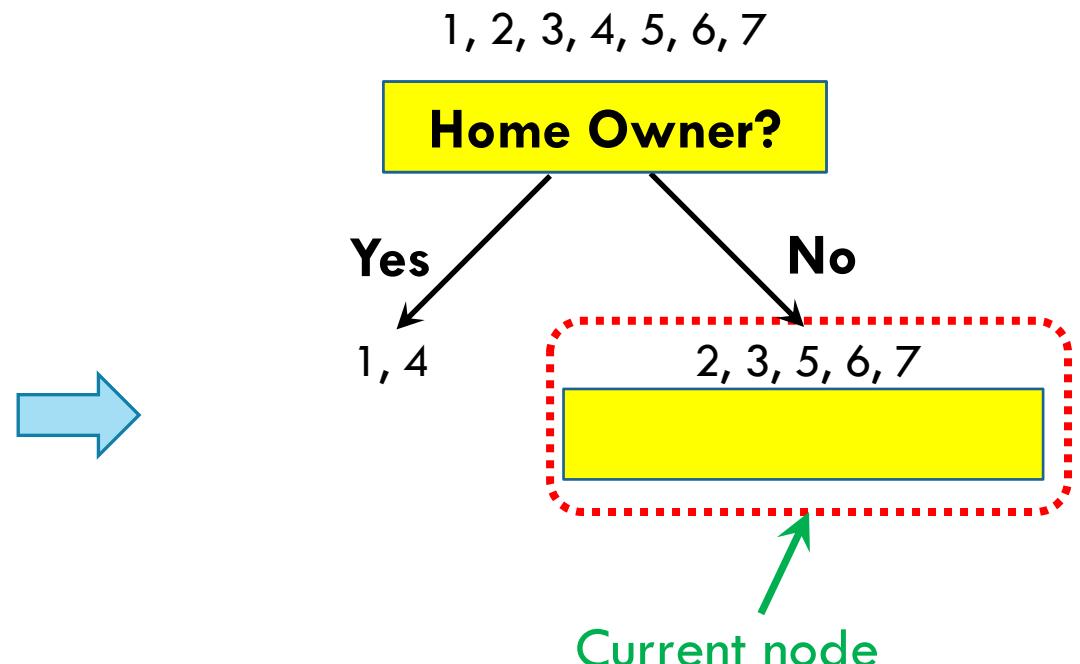


To be defined later

2. Find the “**best quality**” criterion, which partitions the data into 2 subsets

# TRAINING DECISION TREE

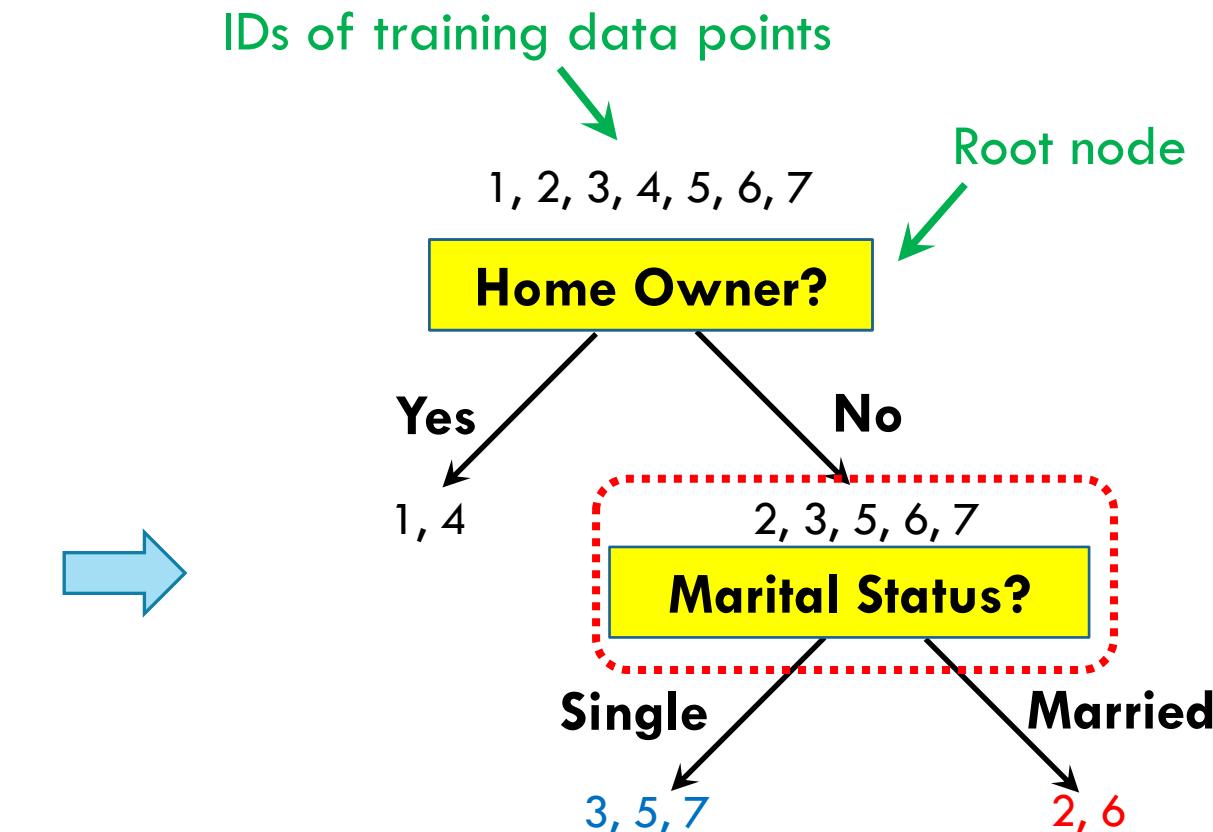
ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
2	No	Married	100K	No
3	No	Single	170K	No
5	No	Single	75K	Yes
6	No	Married	160K	No
7	No	Single	50K	Yes



3. Continue recursively onto the resulting subsets (i.e. splitting them further)

# TRAINING DECISION TREE

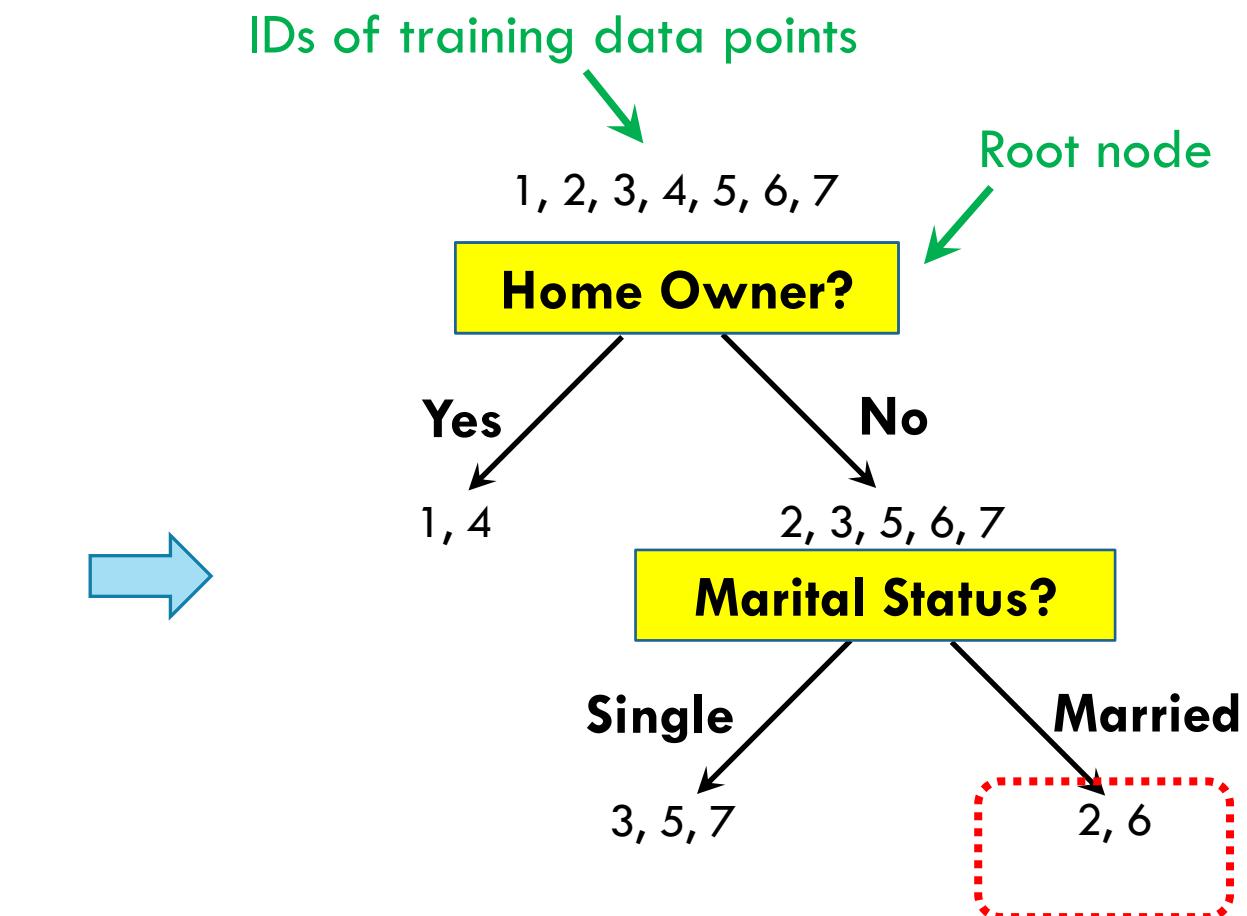
ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
2	No	Married	100K	No
3	No	Single	170K	No
5	No	Single	75K	Yes
6	No	Married	160K	No
7	No	Single	50K	Yes



3. Continue recursively onto the resulting subsets (i.e. splitting them further)

# TRAINING DECISION TREE

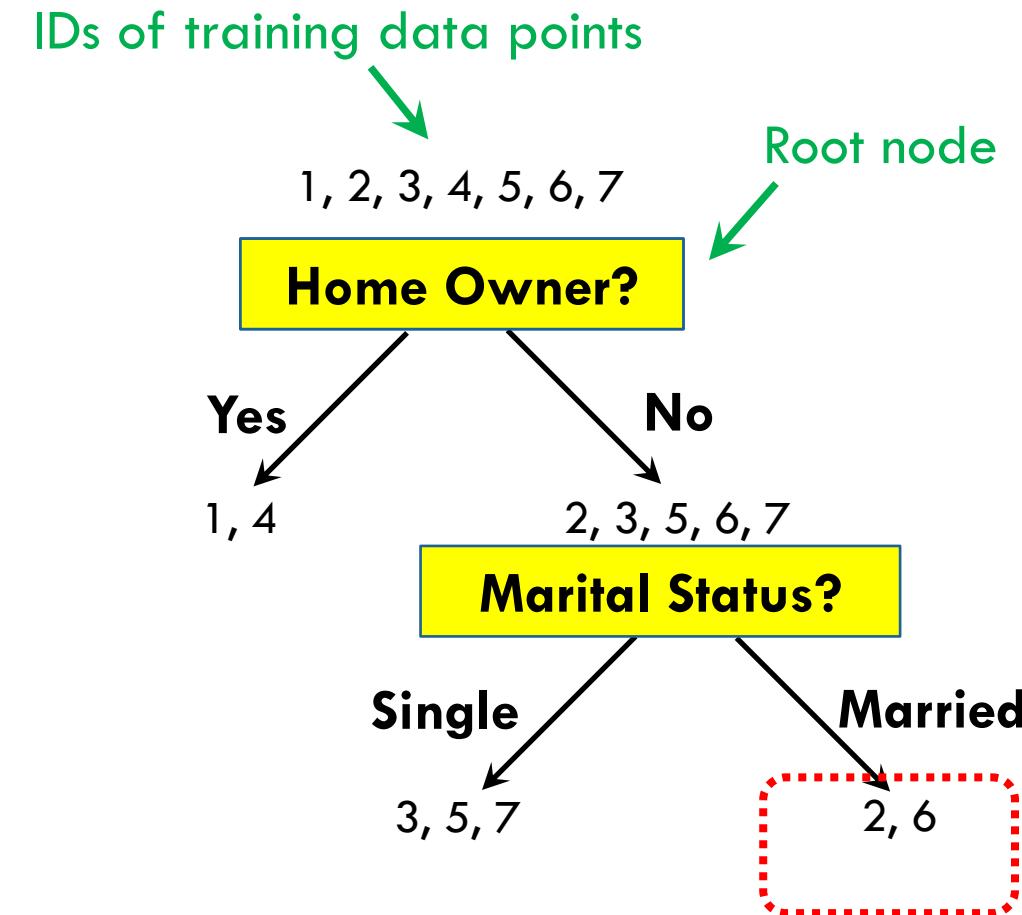
ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
2	No	Married	100K	No
6	No	Married	160K	No



3. Continue recursively onto the resulting subsets (i.e. splitting them further)

# TRAINING DECISION TREE

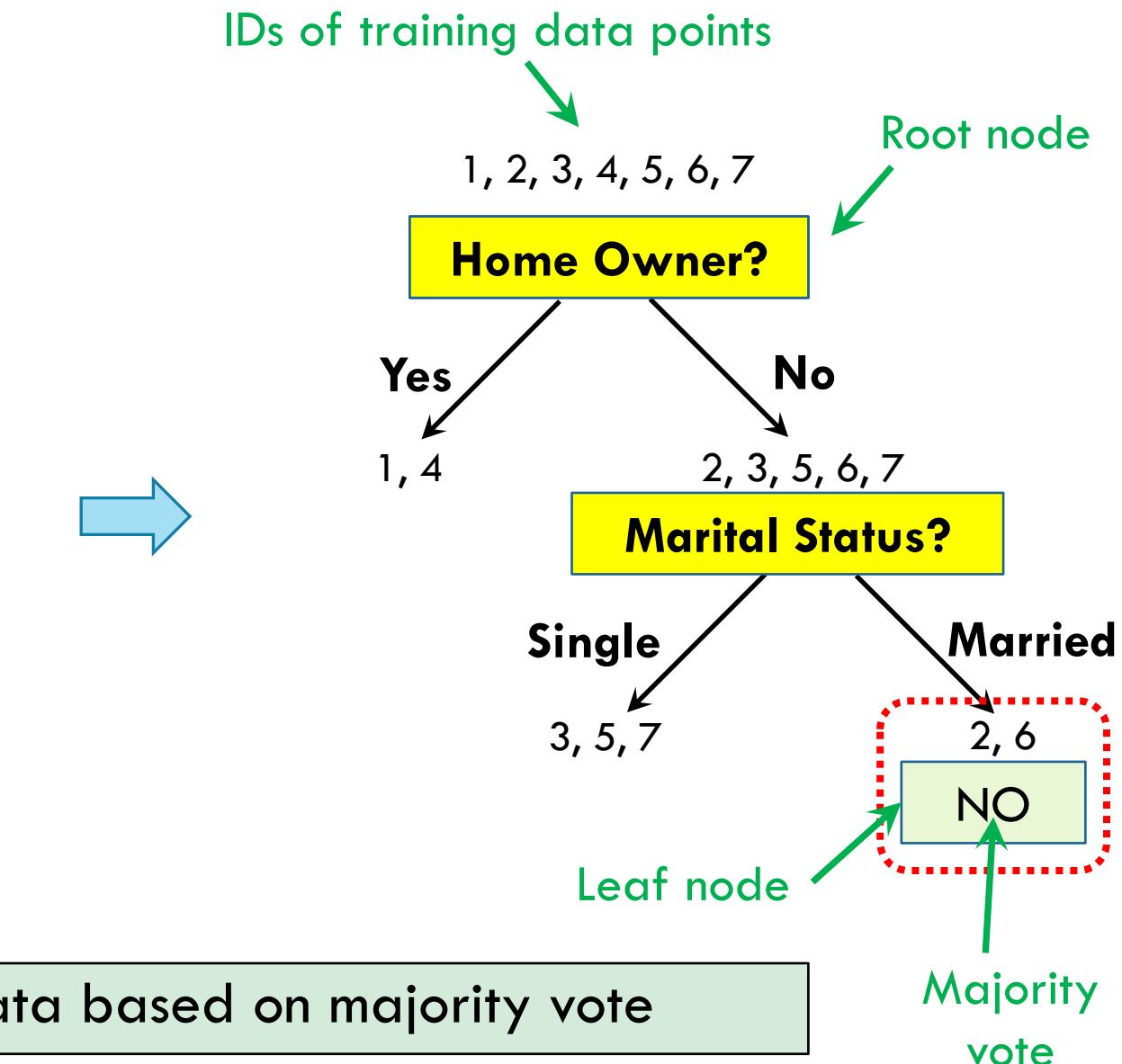
ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
2	No	Married	100K	No
6	No	Married	160K	No



4. Continue until all points have the same response value, or the tree has reached a **depth limit**

# TRAINING DECISION TREE

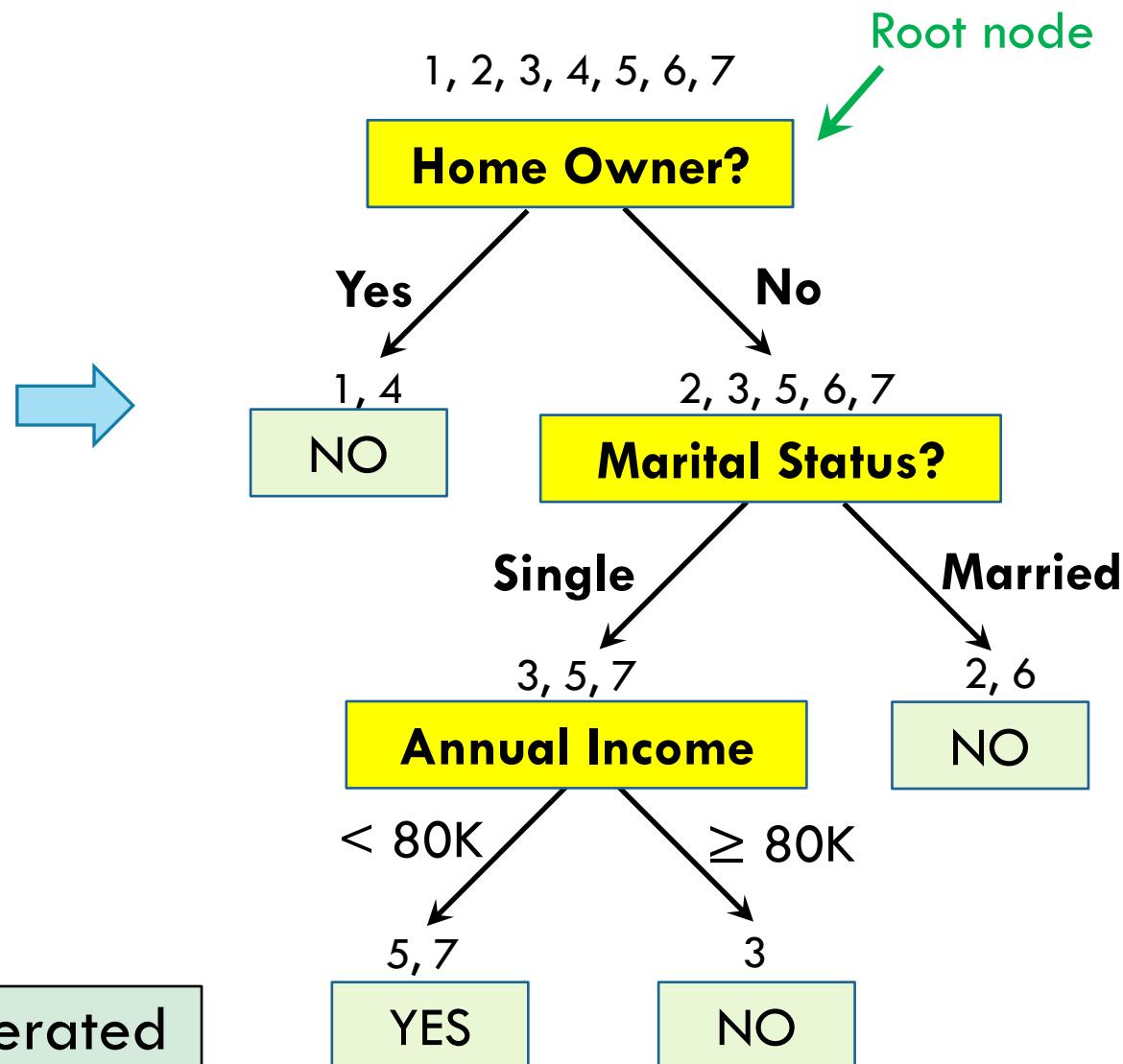
ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
2	No	Married	100K	No
6	No	Married	160K	No



4. On termination, leaf nodes classify the data based on majority vote

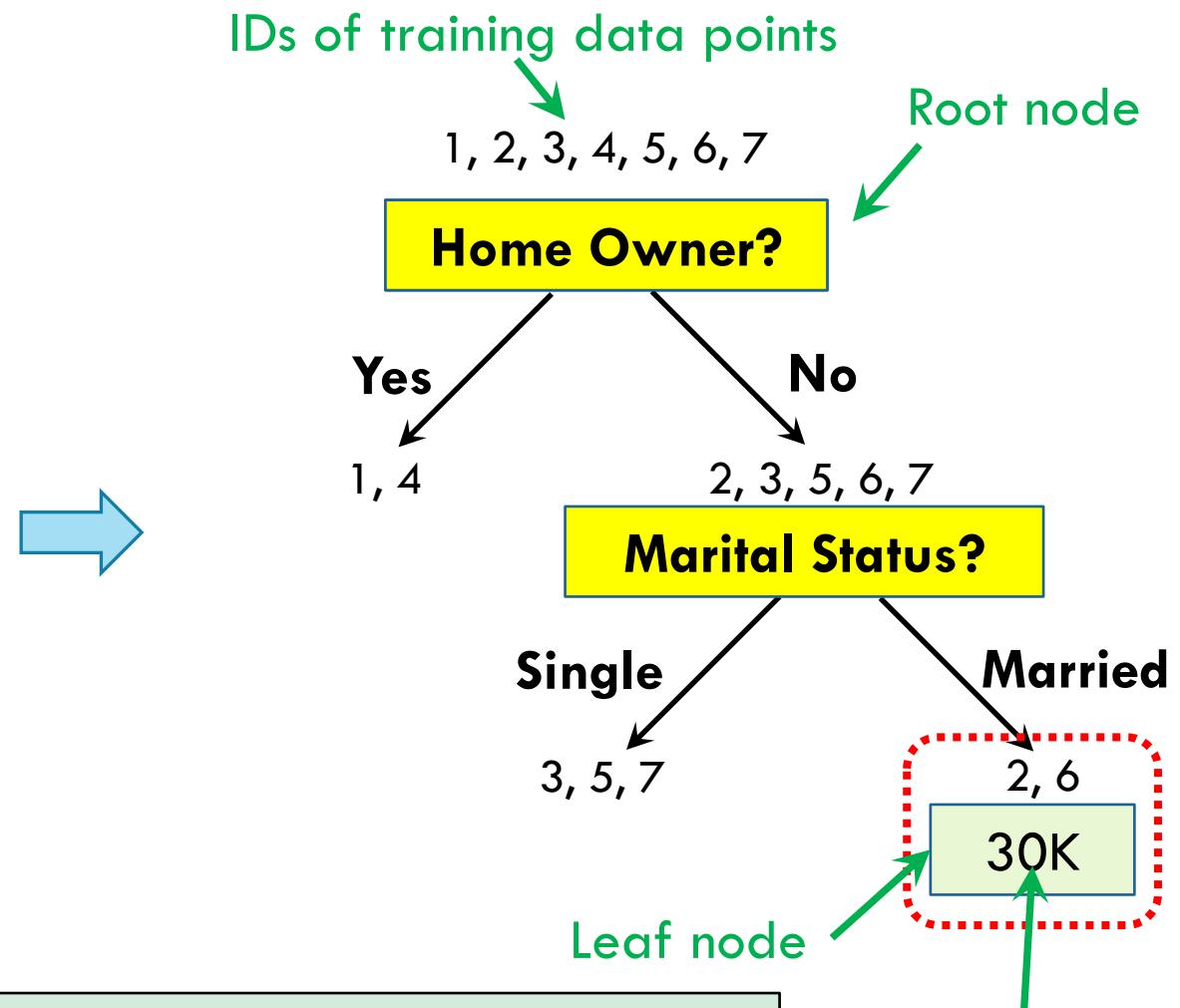
# TRAINING DECISION TREE

ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	170K	No
4	Yes	Married	120K	No
5	No	Single	75K	Yes
6	No	Married	160K	No
7	No	Single	50K	Yes



# TRAINING REGRESSION TREES (FOR NUMERIC RESPONSE VARIABLE)

Numeric response variable				
ID	Home Owner	Marital Status	Annual Income	Defaulted Amount
2	No	Married	100K	20K
6	No	Married	160K	40K



**Regression trees:** instead of using majority vote at leaf nodes, we predict the mean of the response variable for training points assigned to the leaf nodes

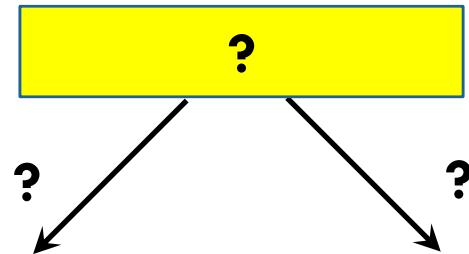
# HOW TO FIND THE BEST SPLIT?

ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	170K	No
4	Yes	Married	120K	No
5	No	Single	75K	Yes
6	No	Married	160K	No
7	No	Single	50K	Yes

## Steps:

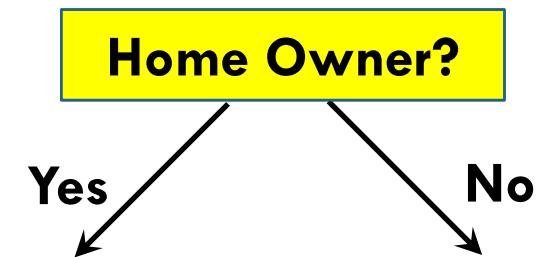
1. Try all **valid splits** of each feature
2. Measure the **quality** of each such split
3. Return the highest quality split

1, 2, 3, 4, 5, 6, 7



# SPLITTING BINARY VARIABLES

ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	170K	No
4	Yes	Married	120K	No
5	No	Single	75K	Yes
6	No	Married	160K	No
7	No	Single	50K	Yes



**Binary Variables:** partition the records into 2 sets based on this variable

# SPLITTING NUMERIC VARIABLES

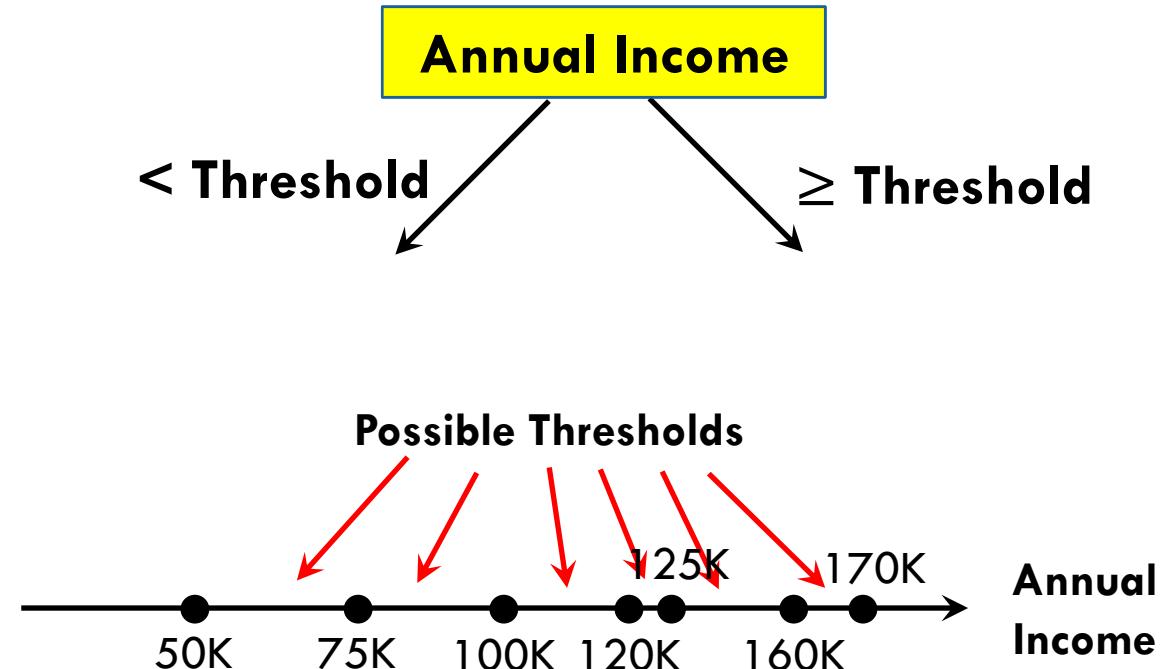
ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	170K	No
4	Yes	Married	120K	No
5	No	Single	75K	Yes
6	No	Married	160K	No
7	No	Single	50K	Yes



**Numeric Variables:** consider all possible partitions into below and above a threshold

# SPLITTING NUMERIC VARIABLES

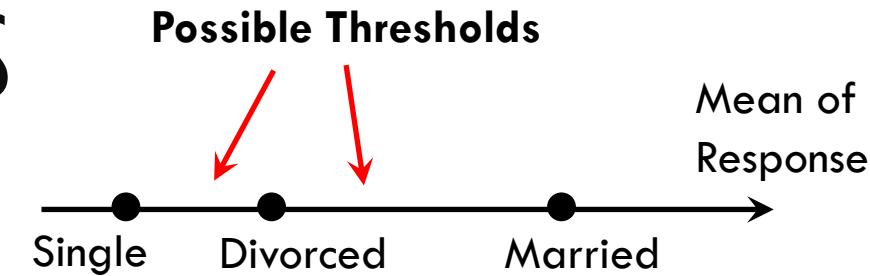
ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	170K	No
4	Yes	Married	120K	No
5	No	Single	75K	Yes
6	No	Married	160K	No
7	No	Single	50K	Yes



**Numeric Variables:** consider all possible partitions into below and above a threshold

# SPLITTING CATEGORICAL VARIABLES

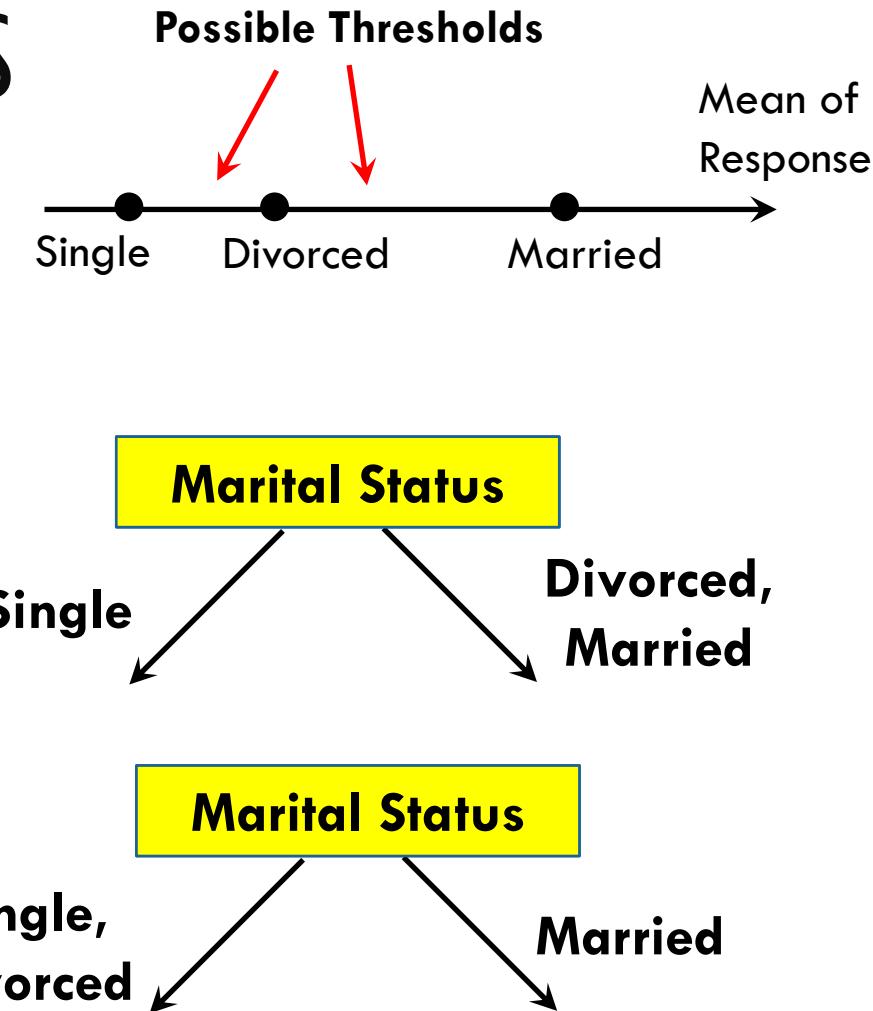
ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	170K	No
4	Yes	Married	120K	No
5	No	Single	75K	Yes
6	No	Married	160K	No
7	No	Single	50K	Yes



**Categorical Variables:** sort the categories by mean response, and treat as continuous variable

# SPLITTING CATEGORICAL VARIABLES

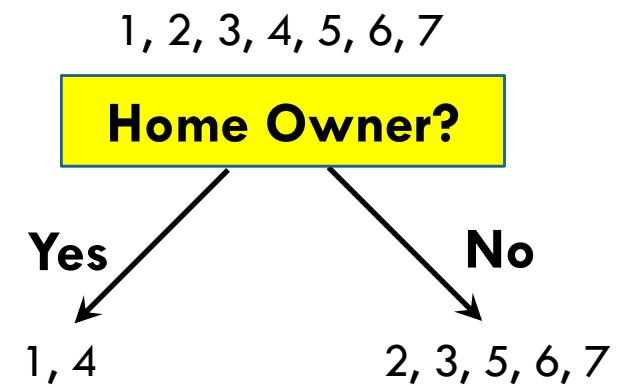
ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	170K	No
4	Yes	Married	120K	No
5	No	Single	75K	Yes
6	No	Married	160K	No
7	No	Single	50K	Yes



**Categorical Variables:** sort the categories by mean response, and treat as continuous variable

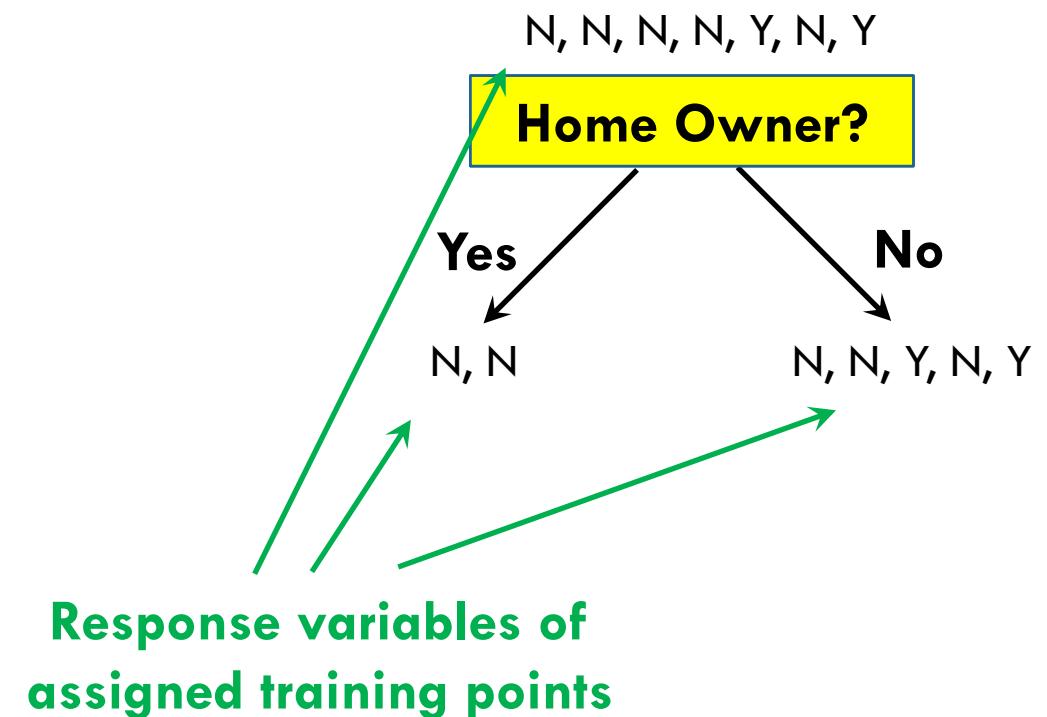
# HOW TO EVALUATE THE QUALITY OF A SPLIT?

ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	170K	No
4	Yes	Married	120K	No
5	No	Single	75K	Yes
6	No	Married	160K	No
7	No	Single	50K	Yes



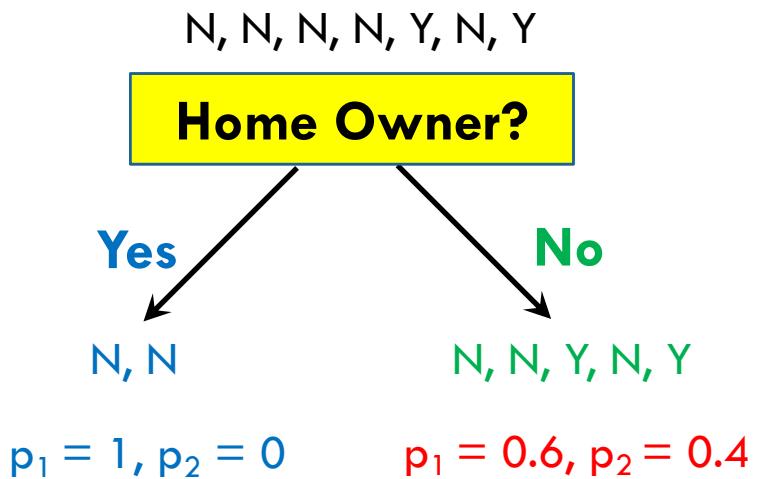
# EVALUATING SPLIT QUALITY: GINI INDEX

ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	170K	No
4	Yes	Married	120K	No
5	No	Single	75K	Yes
6	No	Married	160K	No
7	No	Single	50K	Yes



# EVALUATING SPLIT QUALITY: GINI INDEX

ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	170K	No
4	Yes	Married	120K	No
5	No	Single	75K	Yes
6	No	Married	160K	No
7	No	Single	50K	Yes



$$\text{Gini Index} = 1 - \sum_{i=1}^J p_i^2$$

# EVALUATING SPLIT QUALITY: GINI INDEX

ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	170K	No
4	Yes	Married	120K	No
5	No	Single	75K	Yes
6	No	Married	160K	No
7	No	Single	50K	Yes

N, N, N, N, Y, N, Y

**Home Owner?**

Yes

N, N

No

N, N, Y, N, Y

$$p_1 = 1, p_2 = 0$$

$$\begin{aligned} \text{Gini} &= 1 - 1^2 \\ &= 0 \end{aligned}$$

$$p_1 = 0.6, p_2 = 0.4$$

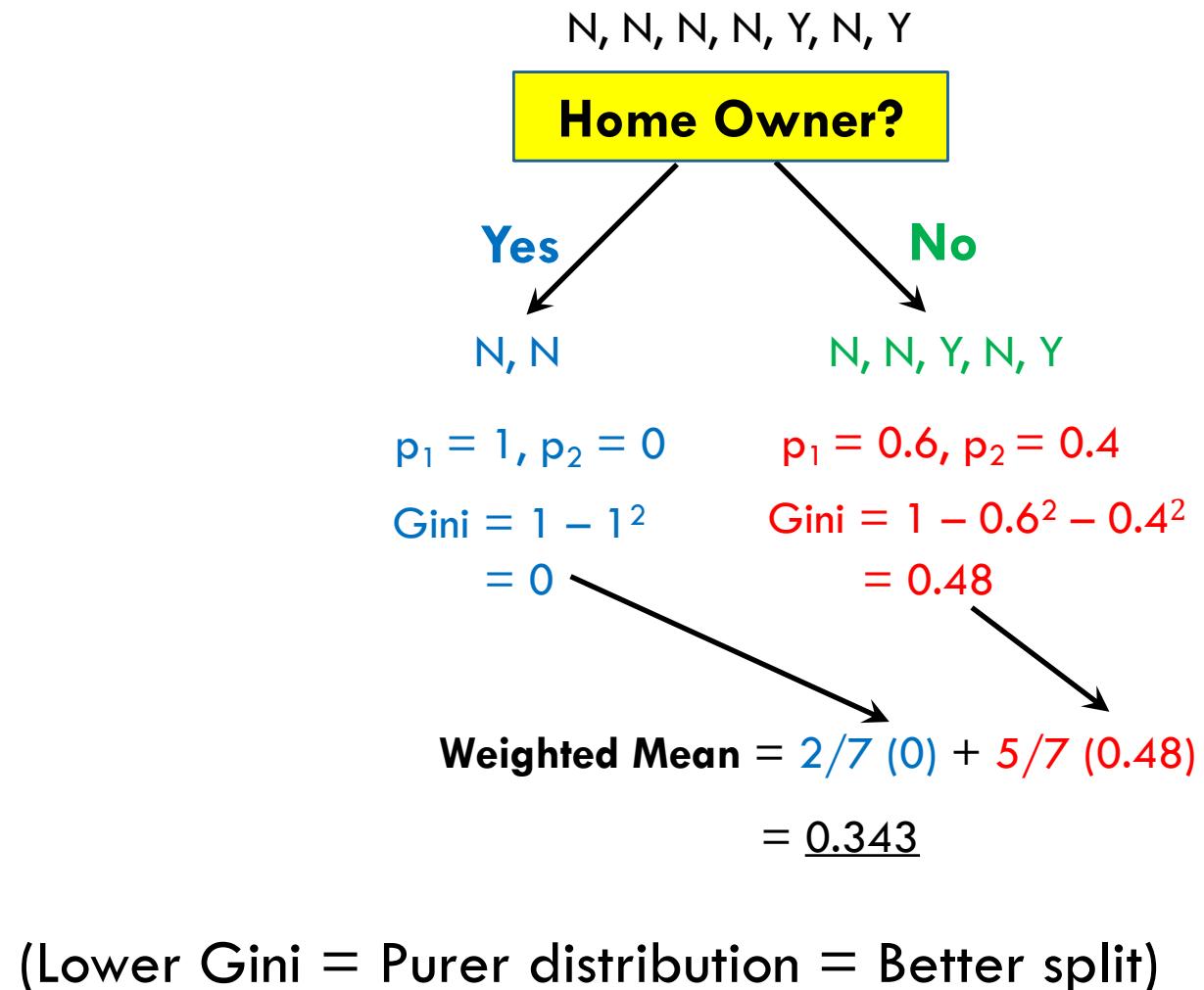
$$\begin{aligned} \text{Gini} &= 1 - 0.6^2 - 0.4^2 \\ &= 0.48 \end{aligned}$$

$$\text{Gini Index} = 1 - \sum_{i=1}^J p_i^2$$

# EVALUATING SPLIT QUALITY: GINI INDEX

ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	170K	No
4	Yes	Married	120K	No
5	No	Single	75K	Yes
6	No	Married	160K	No
7	No	Single	50K	Yes

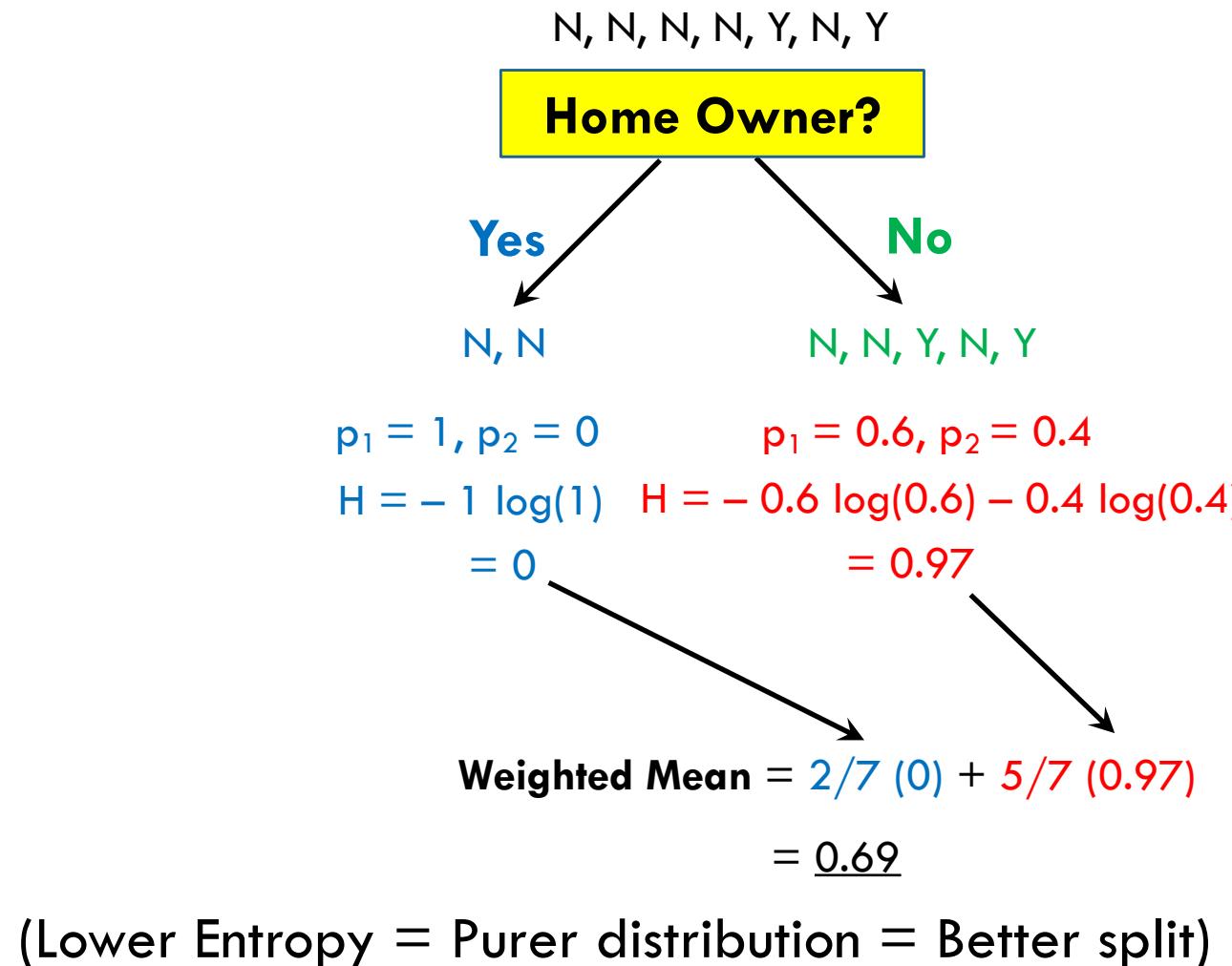
$$\text{Gini Index} = 1 - \sum_{i=1}^J p_i^2$$



# EVALUATING SPLIT QUALITY: ENTROPY

ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	170K	No
4	Yes	Married	120K	No
5	No	Single	75K	Yes
6	No	Married	160K	No
7	No	Single	50K	Yes

$$\text{Entropy: } H = - \sum_{i=1}^J p_i \log_2 p_i$$



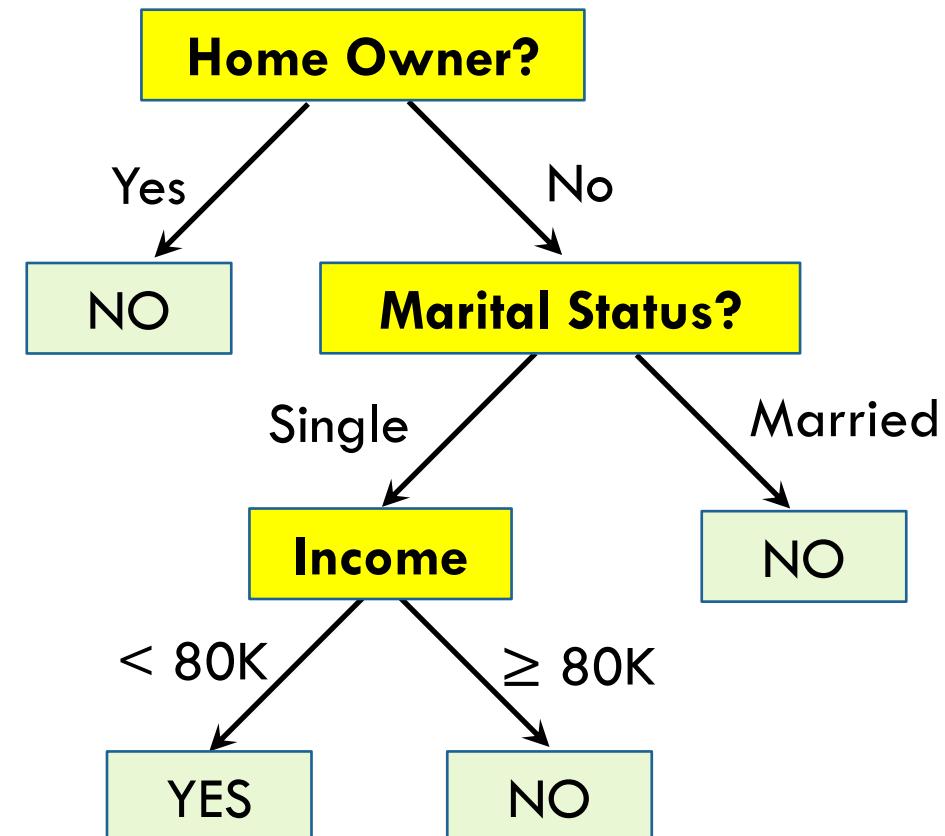
# PROS AND CONS

## Pros

- Fast to train / test
- Interpretable (mimic human reasoning)
- Can handle categorical and continuous features

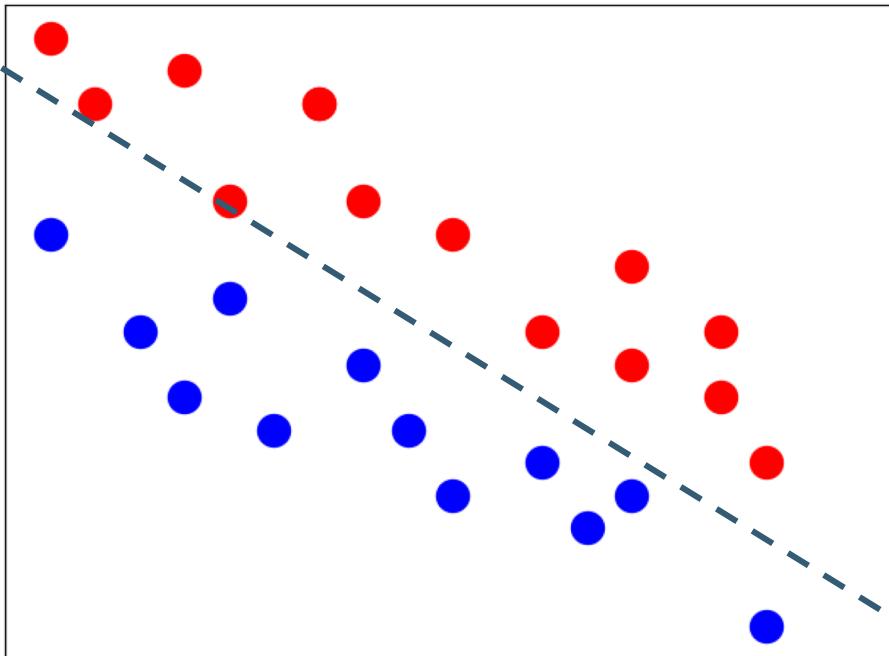
## Cons

- Not robust: small changes in data can result in completely different trees
- Greedy approach may be suboptimal
- Relatively poor accuracy

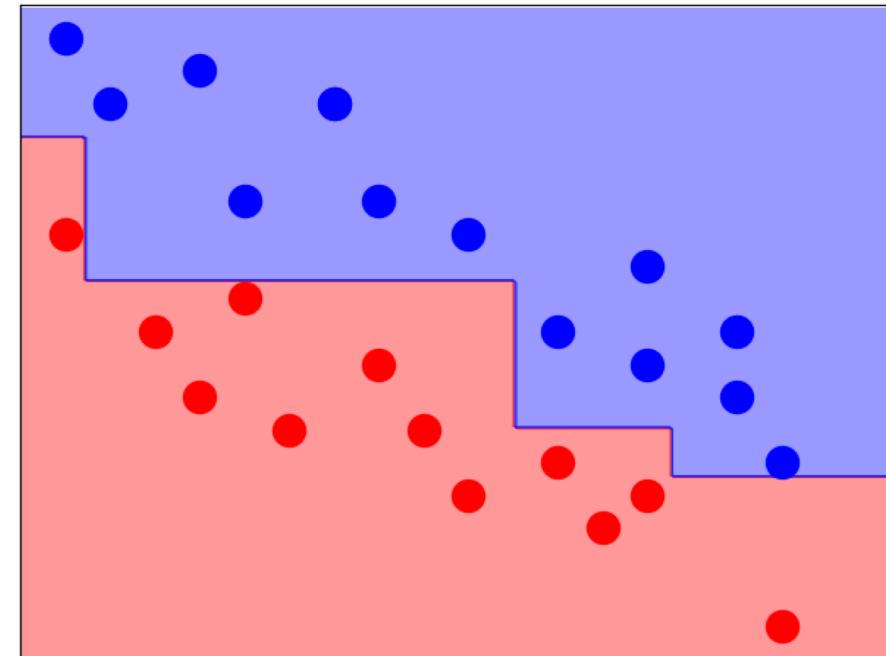


# PROPERTIES – MODELLING INTERACTIONS

Optimal decision boundary

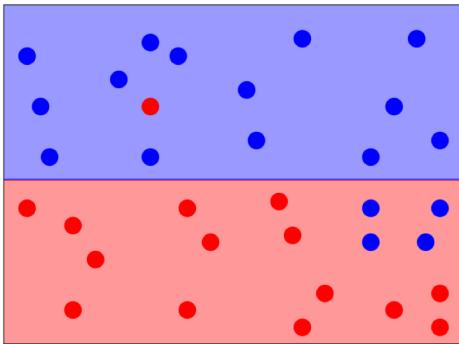


Decision Tree boundaries



# PROPERTIES – UNDERFITTING VS OVERFITTING

Underfitting

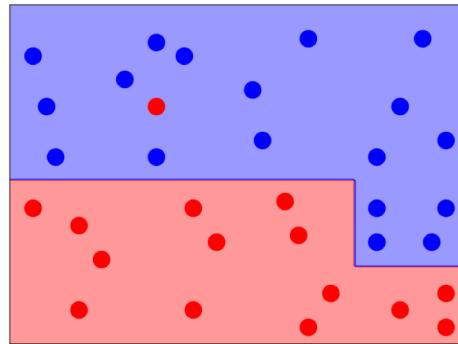


```

X[1] <= 0.485
gini = 0.496
samples = 33
value = [15, 18]
class = y[1]

True      False
gini = 0.346 samples = 18
value = [14, 4] class = y[0]
          gini = 0.124 samples = 15
value = [1, 14] class = y[1]
    
```

Good fit



```

X[1] <= 0.485
gini = 0.496
samples = 33
value = [15, 18]
class = y[1]

True      False
X[0] <= 0.75      X[0] <= 0.35
gini = 0.346      gini = 0.124
samples = 18      samples = 15
value = [14, 4]    value = [1, 14]
class = y[0]       class = y[1]

gini = 0.0
samples = 11
value = [11, 0]
class = y[0]

X[1] <= 0.225      X[0] <= 0.285
gini = 0.49      gini = 0.245
samples = 7      samples = 8
value = [3, 4]    value = [0, 8]
class = y[1]       class = y[1]

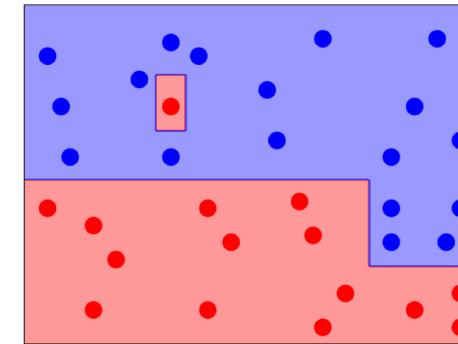
gini = 0.0
samples = 3
value = [3, 0]
class = y[0]

gini = 0.0
samples = 4
value = [0, 4]
class = y[1]

gini = 0.0
samples = 7
value = [1, 6]
class = y[1]

gini = 0.444
samples = 3
value = [1, 1]
class = y[1]
    
```

Overfitting



```

X[1] <= 0.485
gini = 0.496
samples = 33
value = [15, 18]
class = y[1]

True      False
X[0] <= 0.75      X[0] <= 0.35
gini = 0.346      gini = 0.124
samples = 18      samples = 15
value = [11, 0]    value = [1, 14]
class = y[0]       class = y[1]

gini = 0.0
samples = 11
value = [11, 0]
class = y[0]

X[1] <= 0.225      X[0] <= 0.285
gini = 0.49      gini = 0.245
samples = 7      samples = 8
value = [3, 4]    value = [0, 8]
class = y[1]       class = y[1]

gini = 0.0
samples = 3
value = [3, 0]
class = y[0]

gini = 0.0
samples = 4
value = [0, 4]
class = y[1]

gini = 0.0
samples = 7
value = [1, 6]
class = y[1]

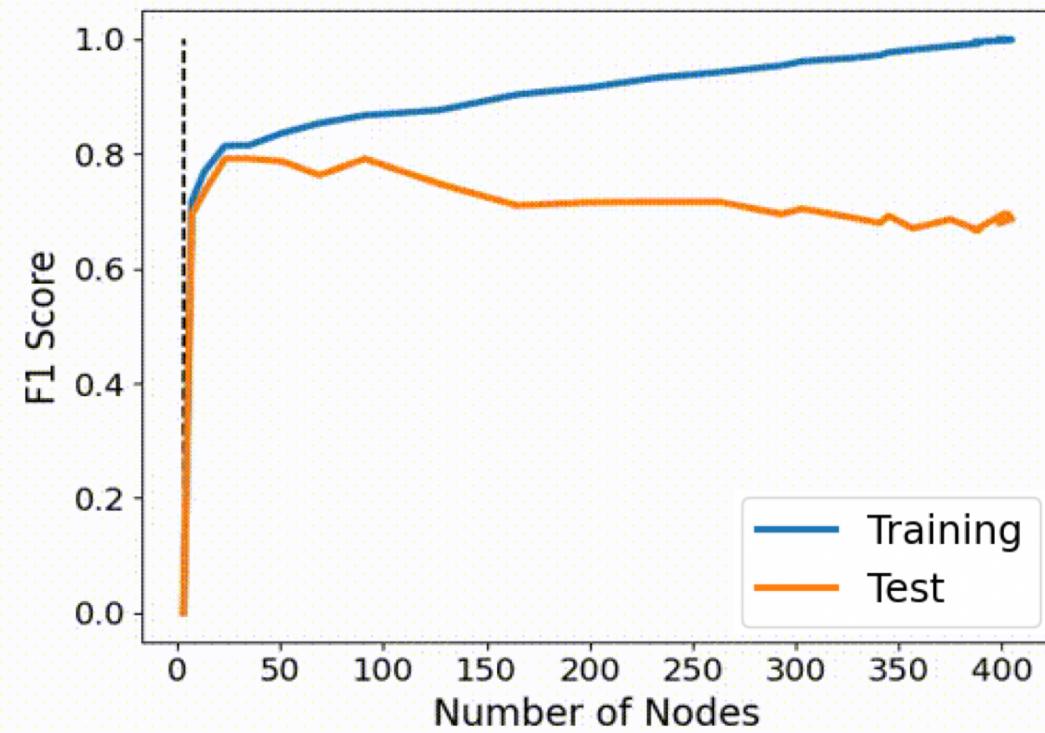
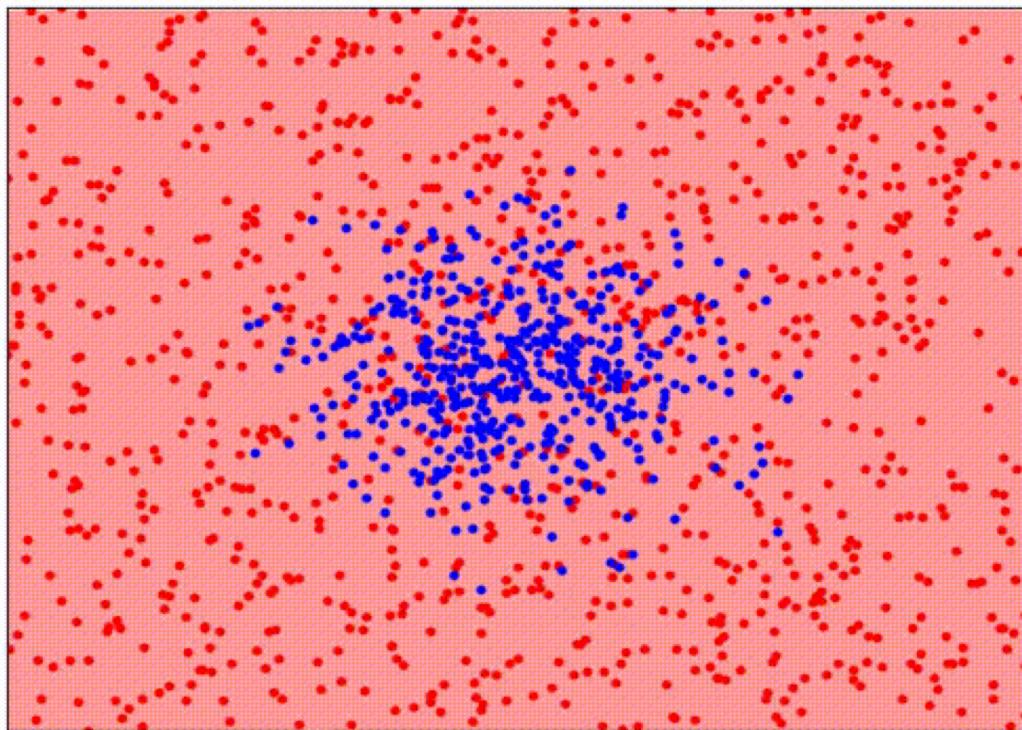
gini = 0.623
samples = 3
value = [1, 2]
class = y[1]

X[1] <= 0.623      X[0] <= 0.795
gini = 0.444      gini = 0.5
samples = 4      samples = 2
value = [0, 4]    value = [1, 1]
class = y[1]       class = y[1]

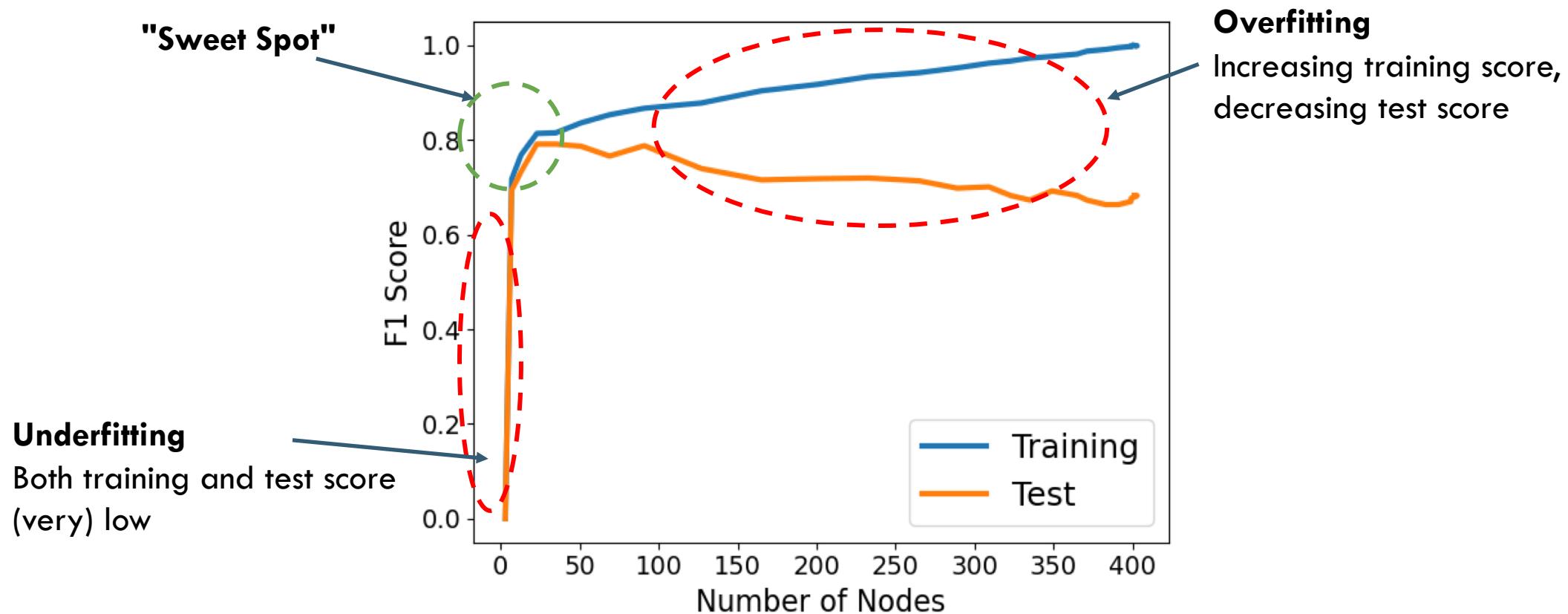
gini = 0.0
samples = 1
value = [1, 0]
class = y[1]

gini = 0.0
samples = 0
value = [0, 1]
class = y[0]
    
```

# PROPERTIES – UNDERFITTING VS OVERFITTING



# PROPERTIES – UNDERFITTING VS OVERFITTING



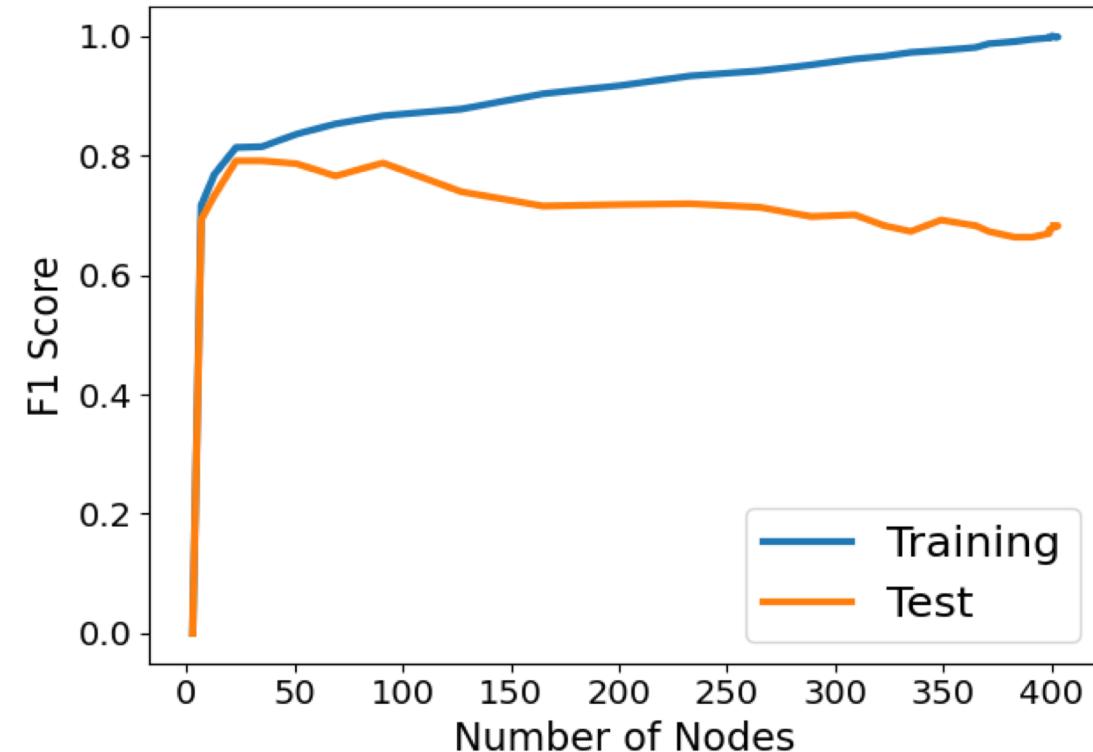
# AVOIDING OVERFITTING

**Issue:** Decision Tree algorithm can always split the training data perfectly (assuming no duplicates)

- It can keep splitting (i.e., increase height of tree) until each leaf contains only one data sample
- One data sample at each leaf → 100% pure

**Solution:** Limit size/height of Decision Tree → Pruning

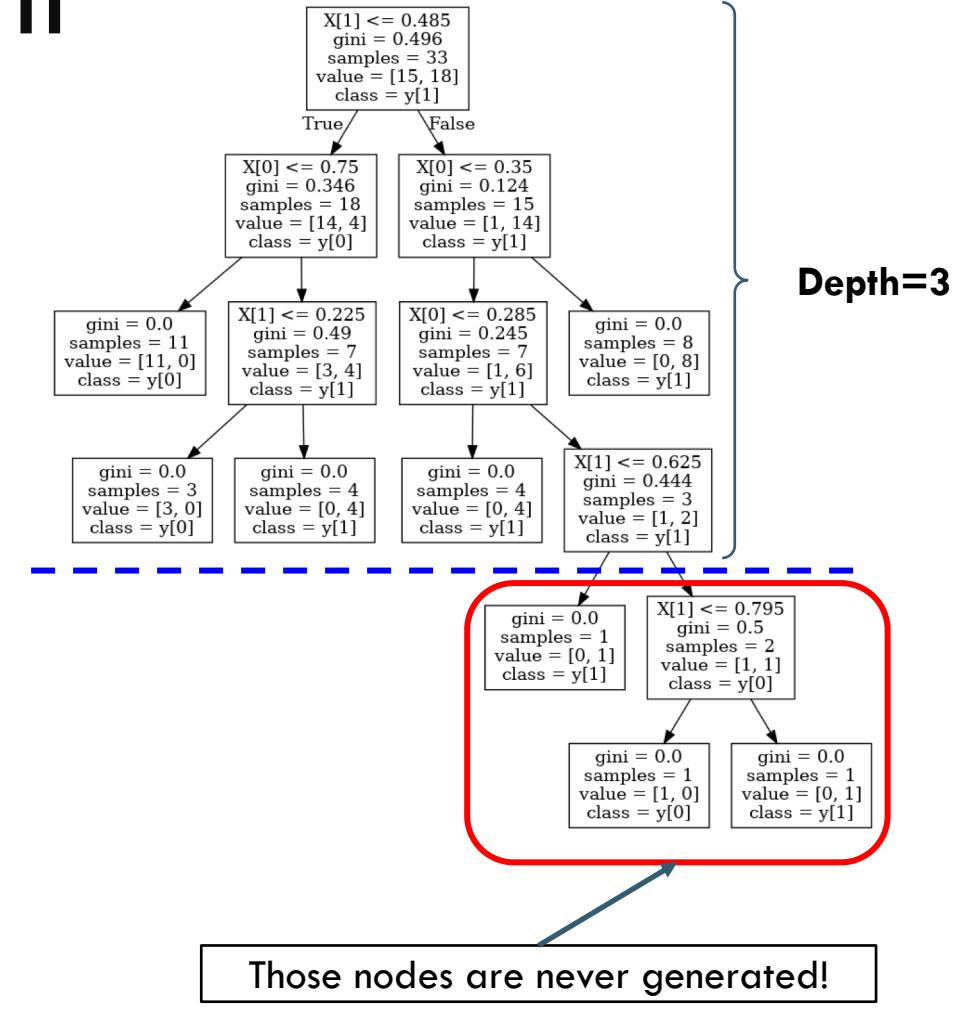
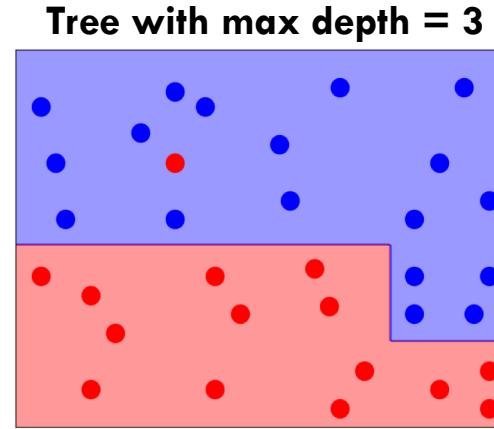
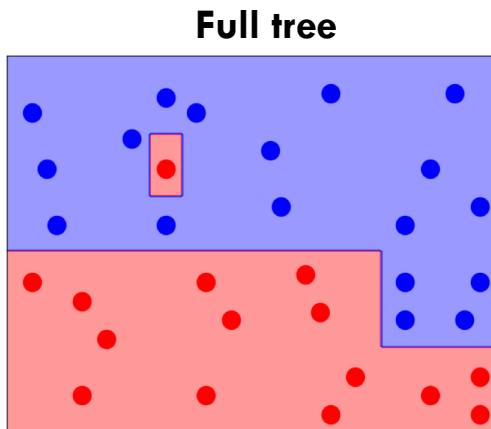
- Pre-pruning: Stop splitting nodes ahead of time
- Post-pruning: Build full tree, but then remove leaves/splits if beneficial
- ... combination of multiple approaches



# PRE-PRUNING: MAXIMUM DEPTH

Stop splitting if maximum depth (set by user) is reached, e.g., depth=3

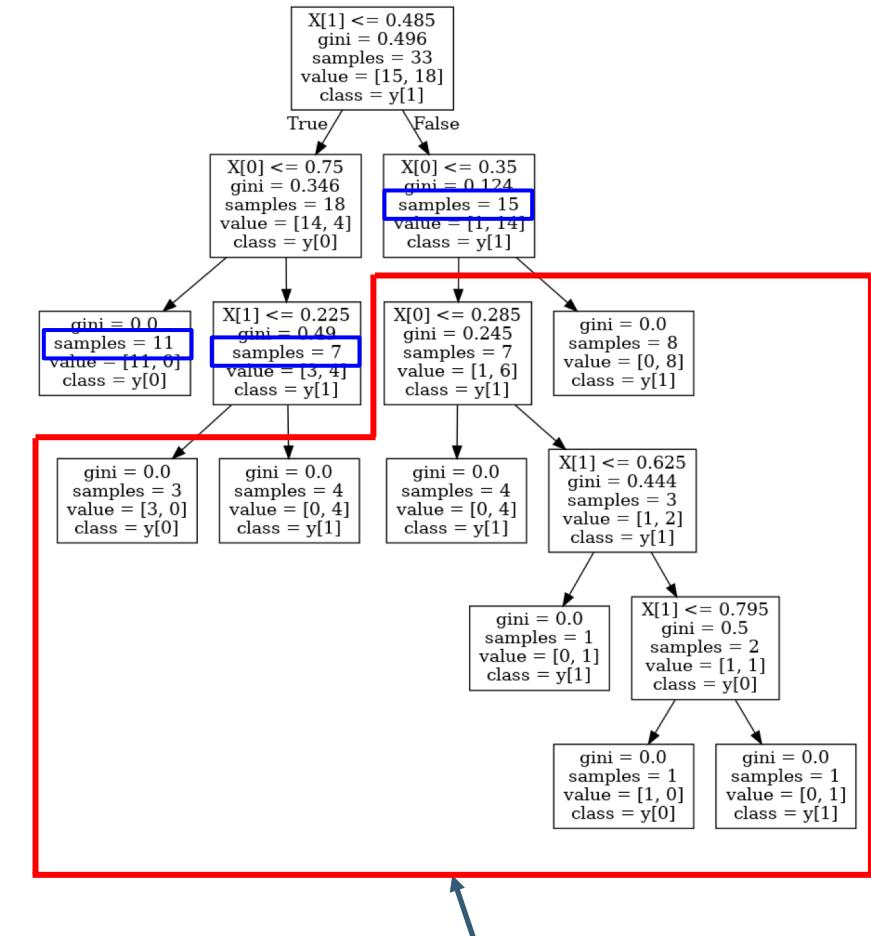
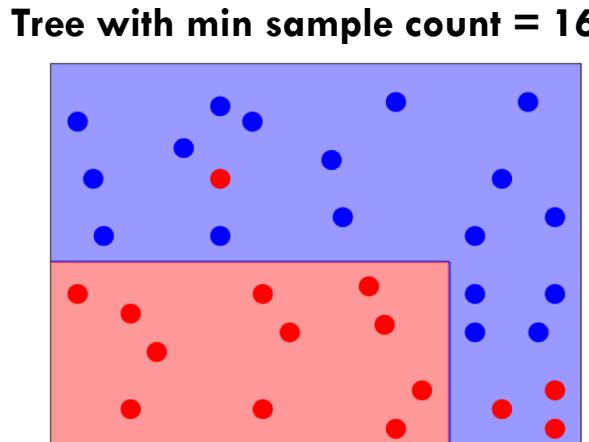
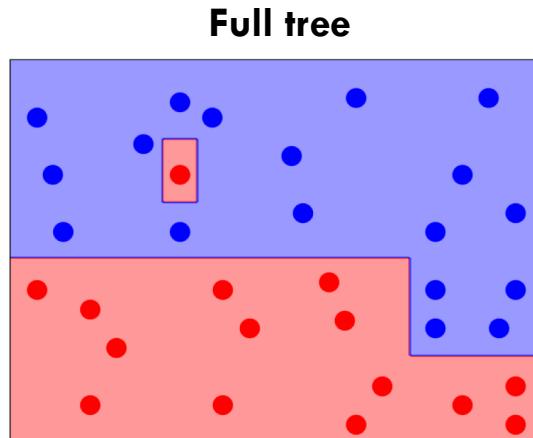
- Corresponds to `max_depth` in scikit-learn



# PRE-PRUNING: MINIMUM SAMPLE COUNT

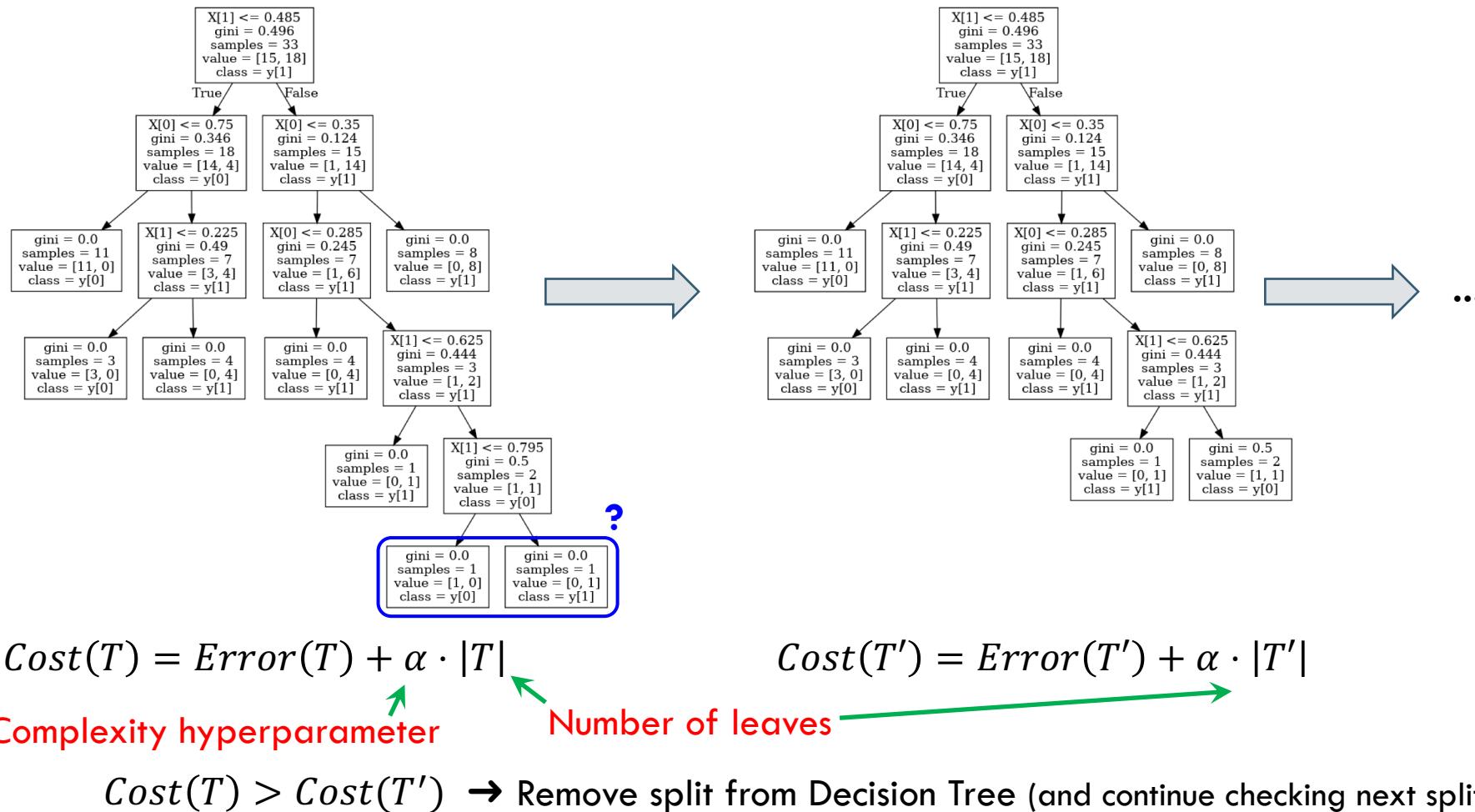
Stop splitting if a node has less than the minimum number of samples: e.g.  
minimum sample count = 16

- Corresponds to `min_samples_split` in scikit-learn



Those nodes are never generated!

# POST-PRUNING: COST-COMPLEXITY PRUNING





# DECISION TREES WHEN RESCALING VARIABLES

**Q:** Is there any change to a trained decision tree's predictions if we shift or rescale any feature (e.g. use annual income in US\$ instead of S\$)?

ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	170K	No
4	Yes	Married	120K	No
5	No	Single	75K	Yes
6	No	Married	160K	No
7	No	Single	50K	Yes

# DECISION TREES WHEN RESCALING VARIABLES



**Q:** Is there any change to a trained decision tree's predictions if we shift or rescale any feature (e.g. use annual income in US\$ instead of S\$)?

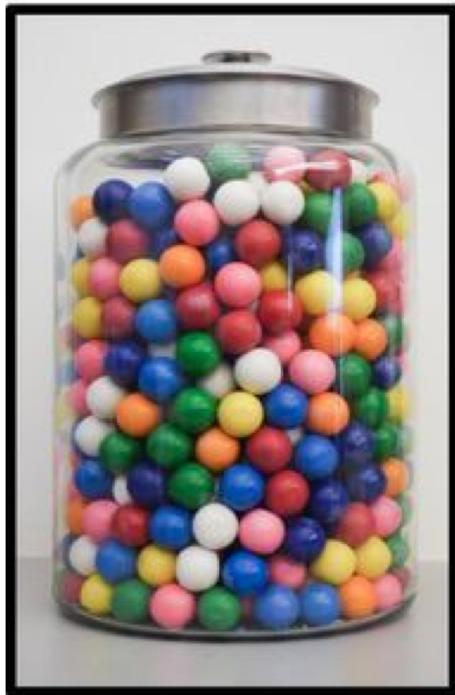
**A:** No, the algorithm does the exact same: e.g. for each cut in US\$, there is a corresponding cut in S\$ that has the same effect. This applies to any order-preserving transformations: e.g. normalization, standardization

ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	170K	No
4	Yes	Married	120K	No
5	No	Single	75K	Yes
6	No	Married	160K	No
7	No	Single	50K	Yes

# CLASSIFICATION OVERVIEW

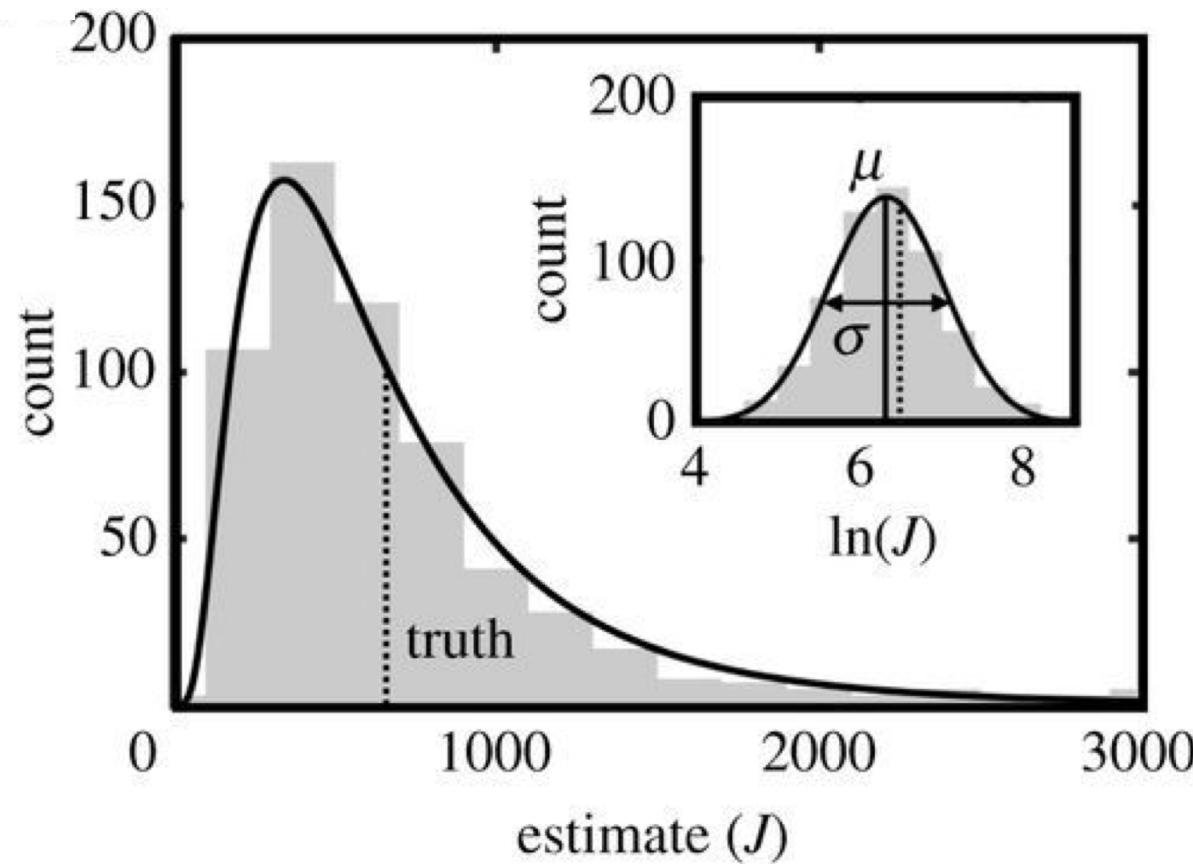
1. Problem Setup
2. Evaluating Classifiers
3. Nearest Neighbor Methods
4. **Trees and Ensembles**
  - a) Decision Trees
  - b) Bagging and Random Forests
  - c) Boosting and Gradient Boosting Machines
5. Logistic Regression
6. Deep Learning

# BAGGING: MOTIVATION ('WISDOM OF THE CROWDS')



$J = 659$  objects

$\ln(J) = 6.5$



# BOOTSTRAP SAMPLING

**Goal:** generate randomly sampled “versions” of the dataset

ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	170K	No
4	Yes	Married	120K	No
5	No	Single	75K	Yes
6	No	Married	160K	No
7	No	Single	50K	Yes

Original Data

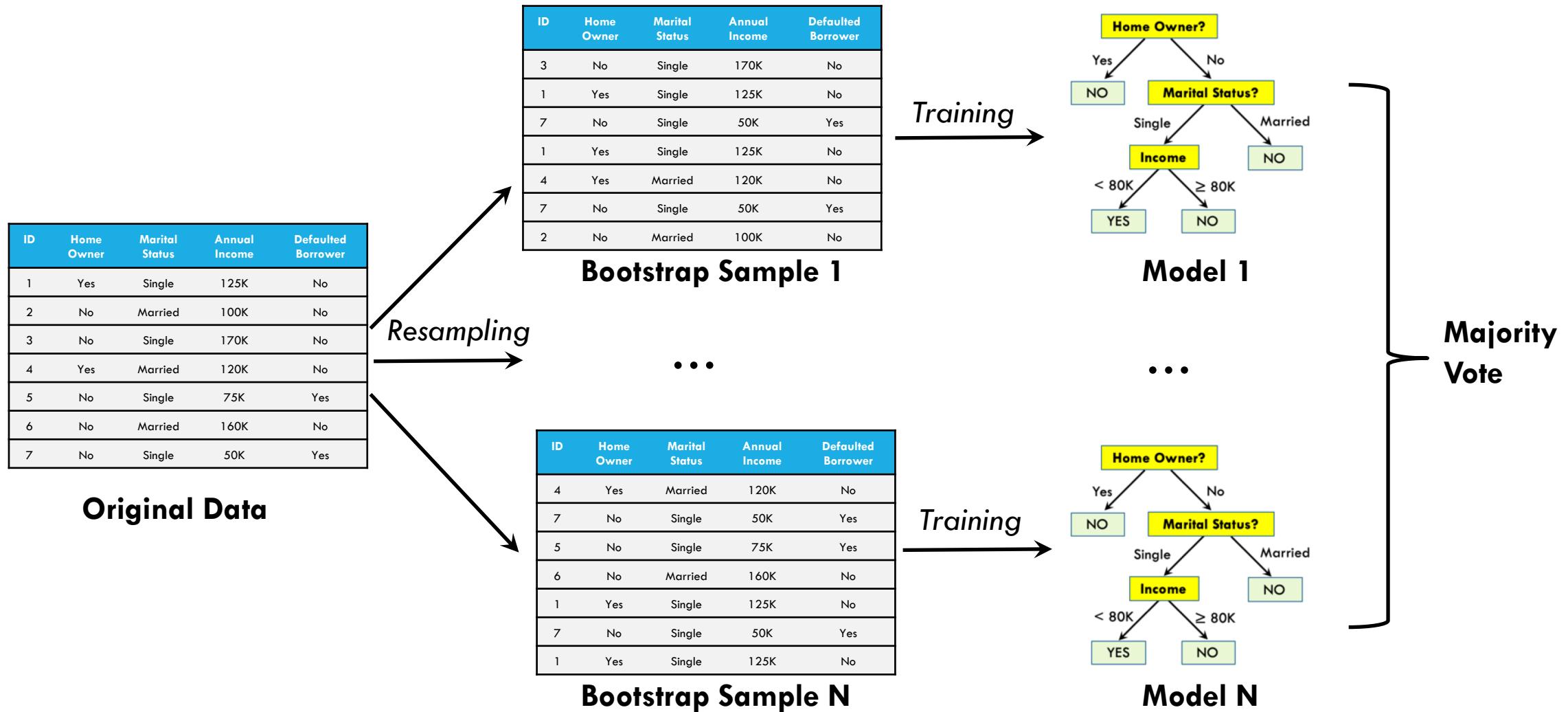


Resample same-sized dataset, with replacement

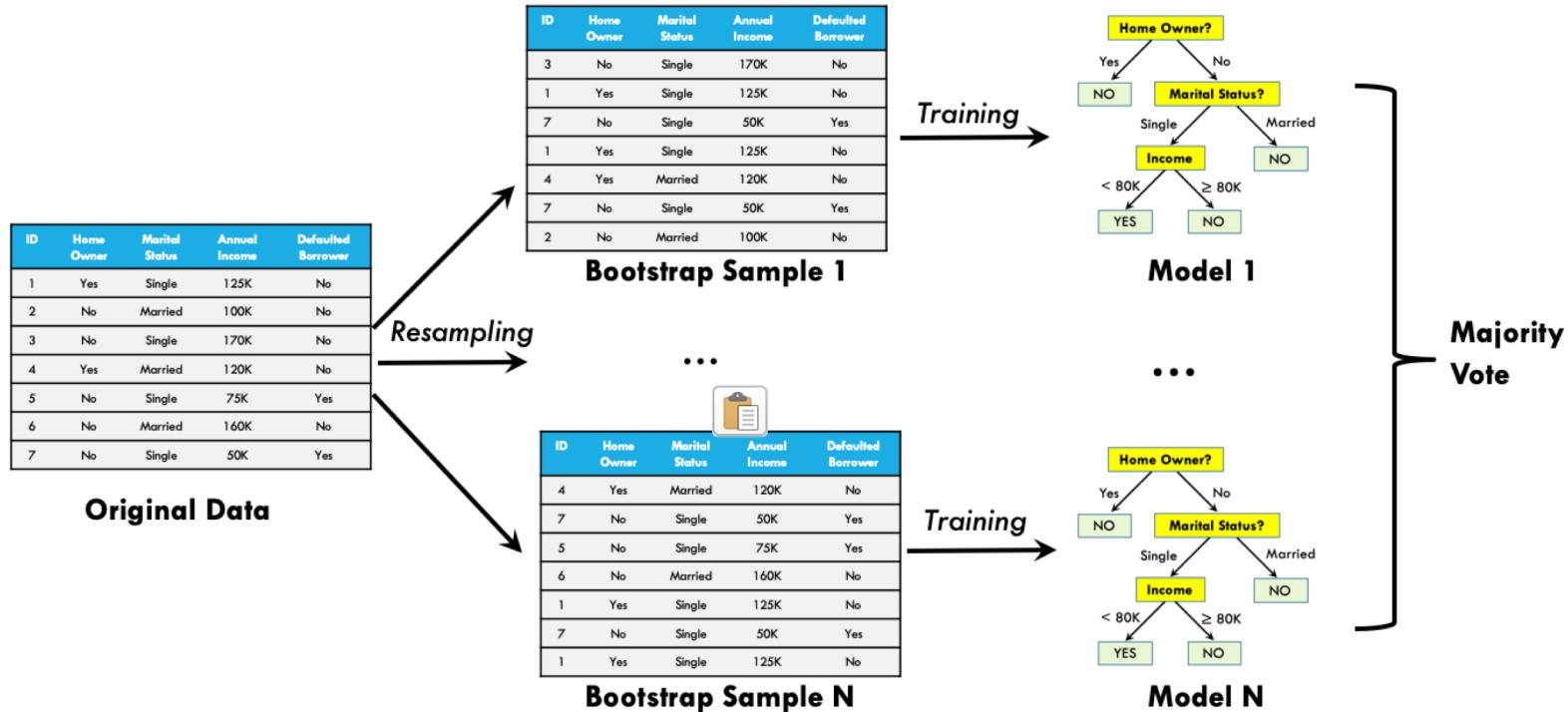
ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
3	No	Single	170K	No
1	Yes	Single	125K	No
7	No	Single	50K	Yes
1	Yes	Single	125K	No
4	Yes	Married	120K	No
7	No	Single	50K	Yes
2	No	Married	100K	No

A Bootstrap Sample

# BOOTSTRAP AGGREGATION ('BAGGING')

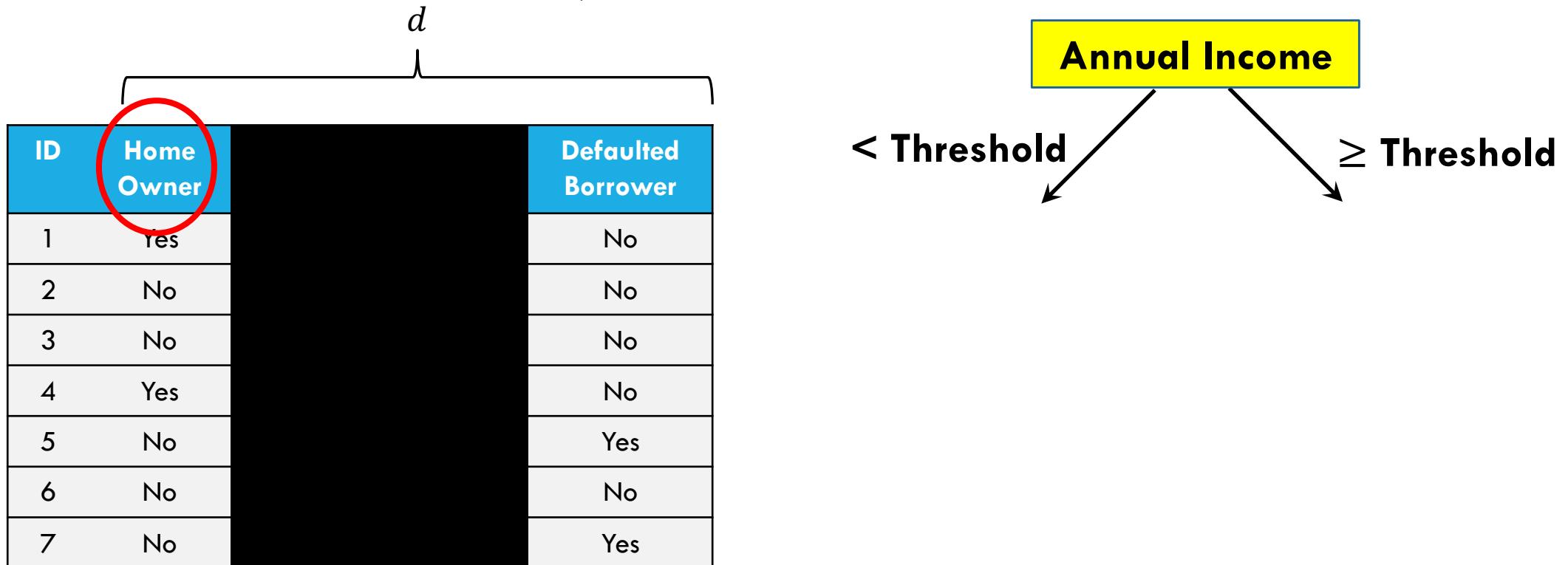


# RANDOM FORESTS: BOOTSTRAP RESAMPLING (“ROW SAMPLING”)



**Bootstrap resampling:** Random forests train an ensemble of decision trees on  $N$  bootstrap samples. (Sometimes, we take bootstrap samples of size smaller than the original dataset size, e.g., based on a user-specified parameter)

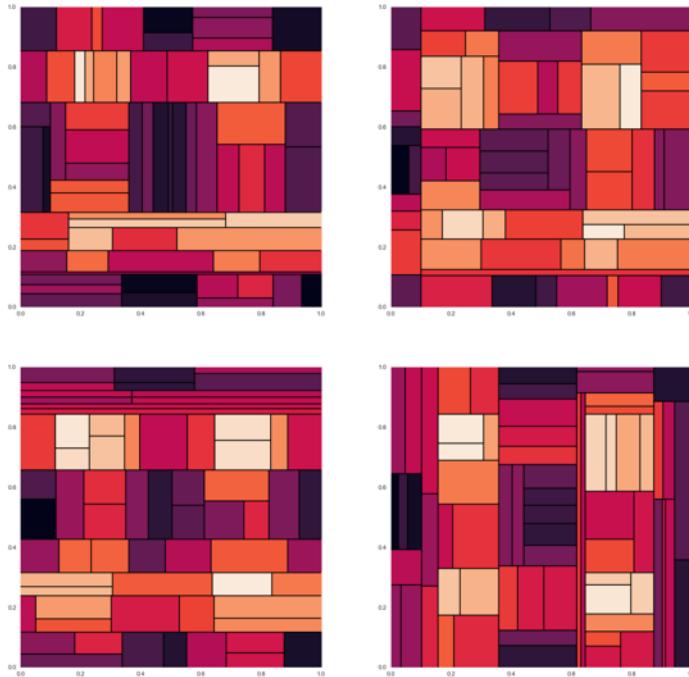
# RANDOM FORESTS: FEATURE SAMPLING (“COLUMN SAMPLING”)



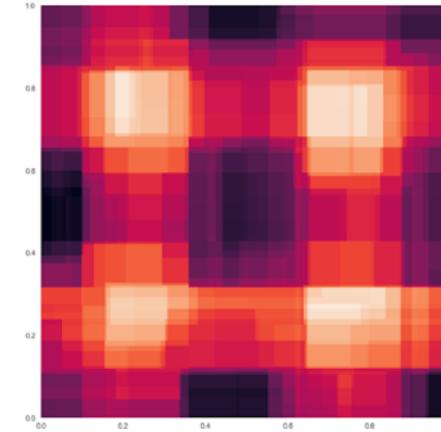
**Feature sampling:** During each split, instead of considering **all  $d$**  features, we only consider a **random subset** of features, usually of size  $\sqrt{d}$

# WHY DO RANDOM FORESTS WORK WELL?

Individual decision tree predictions



Random Forest predictions



Averaging

**Variance Reduction:** individual trees can overfit and give highly variable output. But when averaging them, the predictions are smoother and perform better on test data.

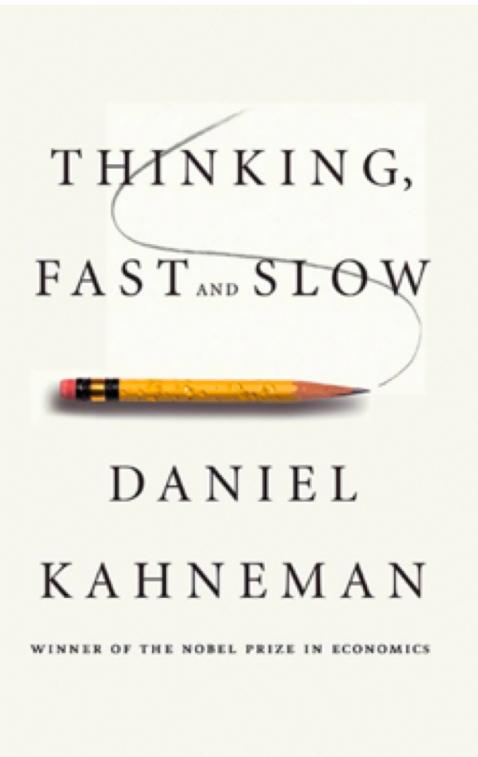
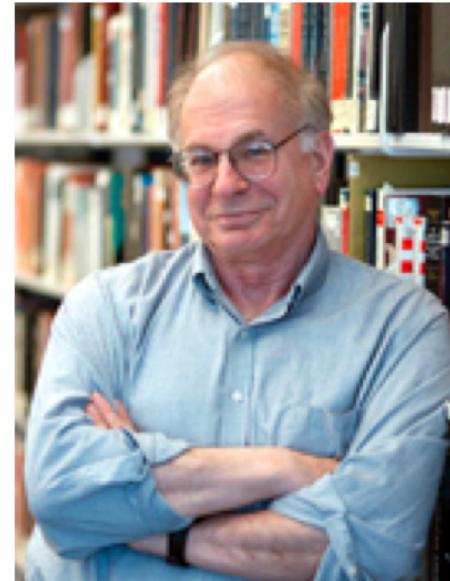
Randomization is important in making the decision trees **decorrelated**

# IMPORTANCE OF DECORRELATION

“To derive the most useful information from multiple sources of evidence, you should always try to make these sources independent of each other.”

“A simple rule can help: before an issue is discussed, all members of the committee should be asked to write a very brief summary of their position. This procedure makes good use of the value of the **diversity of knowledge and opinion in the group**. The standard practice of open discussion gives too much weight to the opinions of those who speak early and assertively, causing others to line up behind them.”

Daniel Kahnemann, *Thinking Fast and Slow*



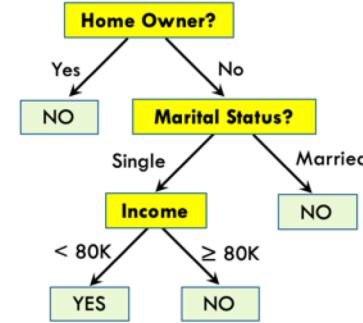
# PROS AND CONS (VS. DECISION TREES)

## Pros

- Variance Reduction: Ensembling many decision trees leads to more stable and accurate predictions
- Accuracy is fairly close to state of the art
- Parallelizable
- Not much tuning required

## Cons

- Less interpretable than decision trees
- Slower than decision trees



**Model 1**

...



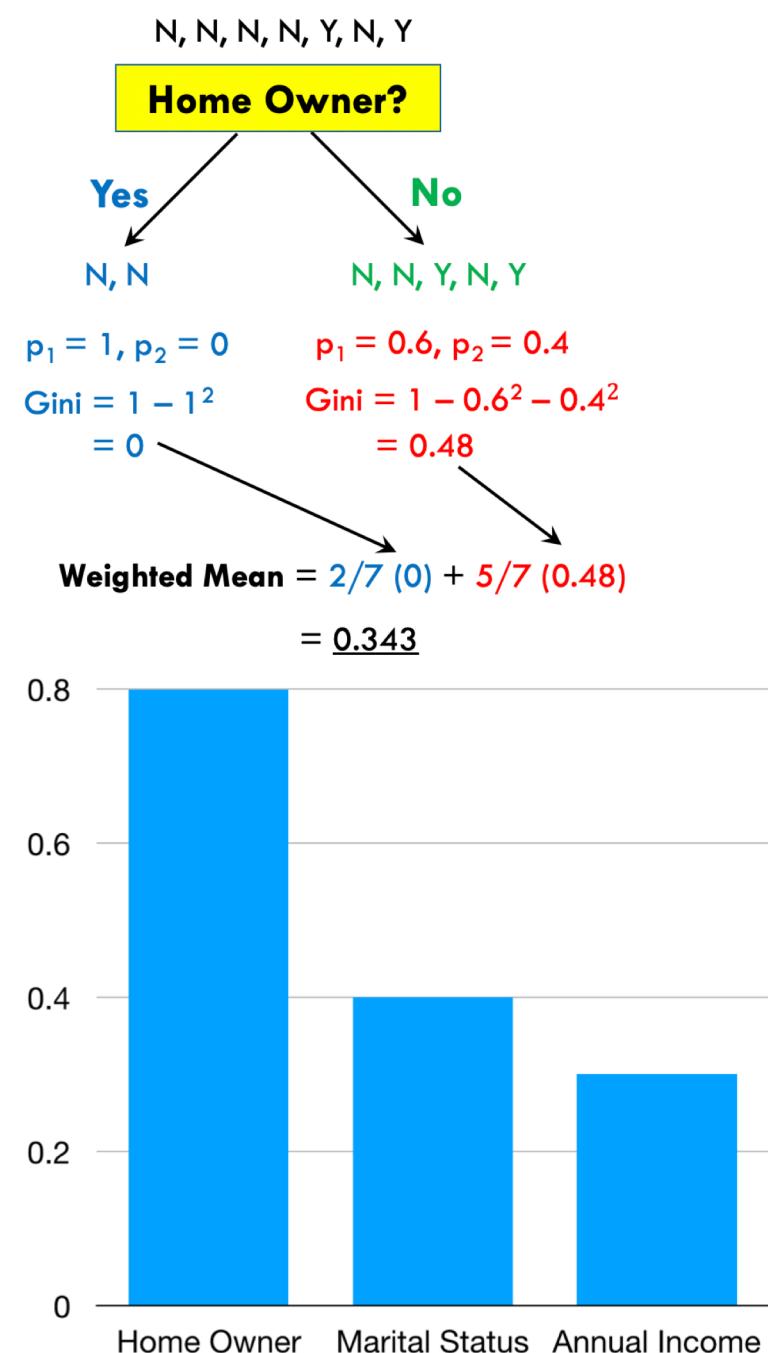
**Model N**

# VARIABLE IMPORTANCE PLOTS

The importance of each variable is measured by the **total reduction in Gini index** brought about by that feature.

More important variables result in greater decrease in Gini index on average.

Computing variable importance is useful for **feature selection**.





# QUIZ: RANDOM FOREST HYPERPARAMETERS

**Q:** Which of the following tends to reduce overfitting?

1. Increasing number of trees
2. Increasing depth of trees
3. Decreasing the number of features considered when selecting each split



**Model 1**

• • •



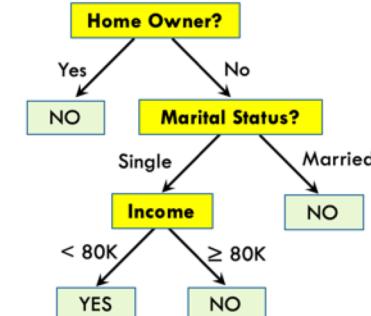
**Model N**



# QUIZ: RANDOM FOREST HYPERPARAMETERS

**Q:** Which of the following tends to reduce overfitting?

1. Increasing number of trees 
2. Increasing depth of trees 
3. Decreasing the number of features considered when selecting each split 



...

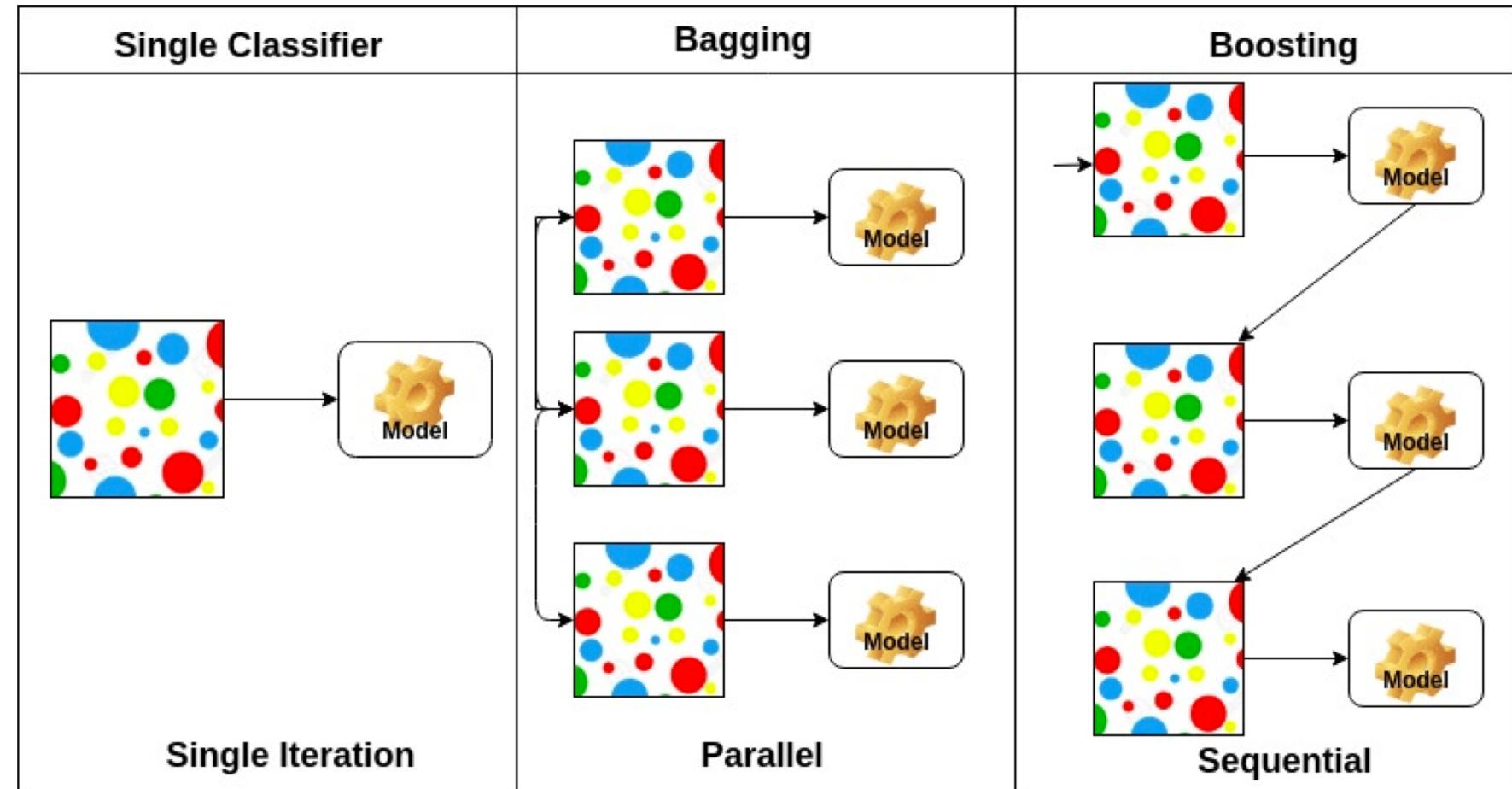


# CLASSIFICATION OVERVIEW

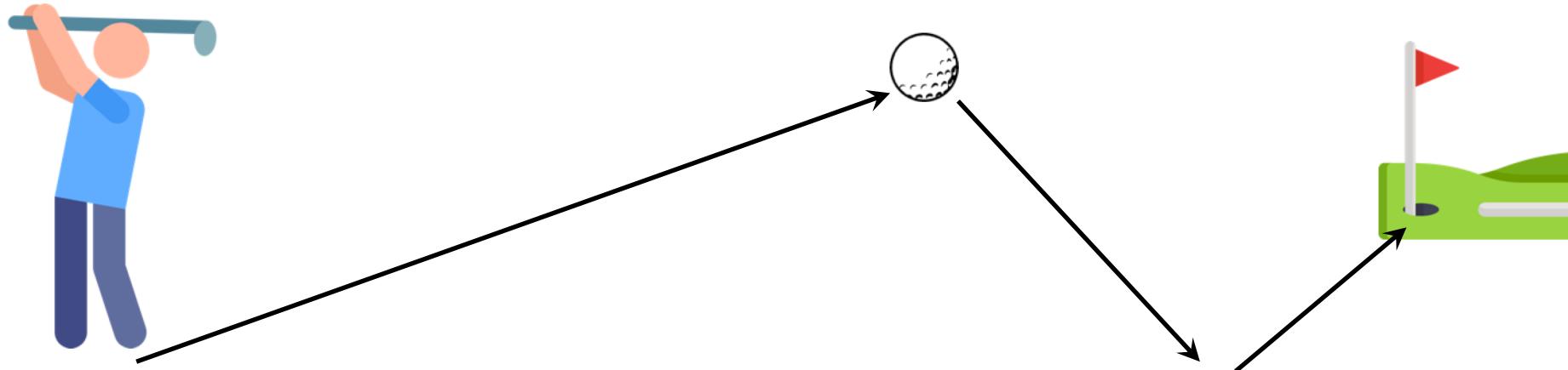
1. Problem Setup
2. Evaluating Classifiers
3. Nearest Neighbor Methods
4. **Trees and Ensembles**
  - a) Decision Trees
  - b) Bagging and Random Forests
  - c) Boosting and Gradient Boosting Machines
5. Logistic Regression
6. Deep Learning

# BAGGING VS BOOSTING

Bagging operates on bootstrap samples **independently**, but boosting is **sequential**: each new classifier “corrects for the mistakes” of the previous classifiers



# BOOSTING: ITERATIVE PROCESS



Each new classifier is trained to “correct for the mistakes” of the previous set of classifiers, gradually getting closer to the ideal classifier

# GRADIENT BOOSTED DECISION TREES

Current Fitted Model

ID	Home Owner	Marital Status	Annual Income	y	r <sub>0</sub>
1	Yes	Single	125K	0	0
2	No	Married	100K	0.5	0.5
3	No	Single	170K	0.3	0.3
4	Yes	Married	120K	1	1
5	No	Single	75K	0.9	0.9
6	No	Married	160K	0.4	0.4
7	No	Single	50K	0.3	0.3

Original Data

Residual  
(Initial)

1. Initial residuals are the original response data

# GRADIENT BOOSTED DECISION TREES

Current Fitted Model

ID	Home Owner	Marital Status	Annual Income	y
1	Yes	Single	125K	0
2	No	Married	100K	0.5
3	No	Single	170K	0.3
4	Yes	Married	120K	1
5	No	Single	75K	0.9
6	No	Married	160K	0.4
7	No	Single	50K	0.3

Original Data (X)

	$r_0$	$f_1$
	0	0.2
	0.5	0.3
	0.3	0.3
	1	1
	0.9	1
	0.4	0.4
	0.3	0.3

Residual  
(Initial)

Fit tree with  
response  $r_0$

2. Fit a decision tree to predict residuals  $r_0$

# GRADIENT BOOSTED DECISION TREES

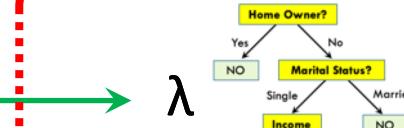
ID	Home Owner	Marital Status	Annual Income	y
1	Yes	Single	125K	0
2	No	Married	100K	0.5
3	No	Single	170K	0.3
4	Yes	Married	120K	1
5	No	Single	75K	0.9
6	No	Married	160K	0.4
7	No	Single	50K	0.3

Original Data (X)

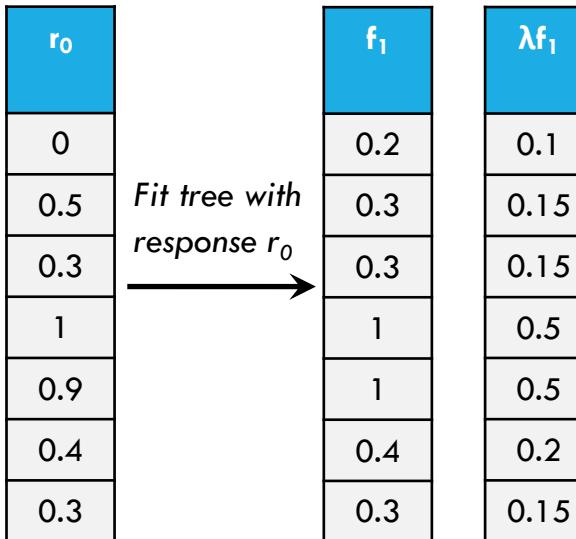
Residual  
(Initial)

Learning rate

$$\lambda = 0.5$$

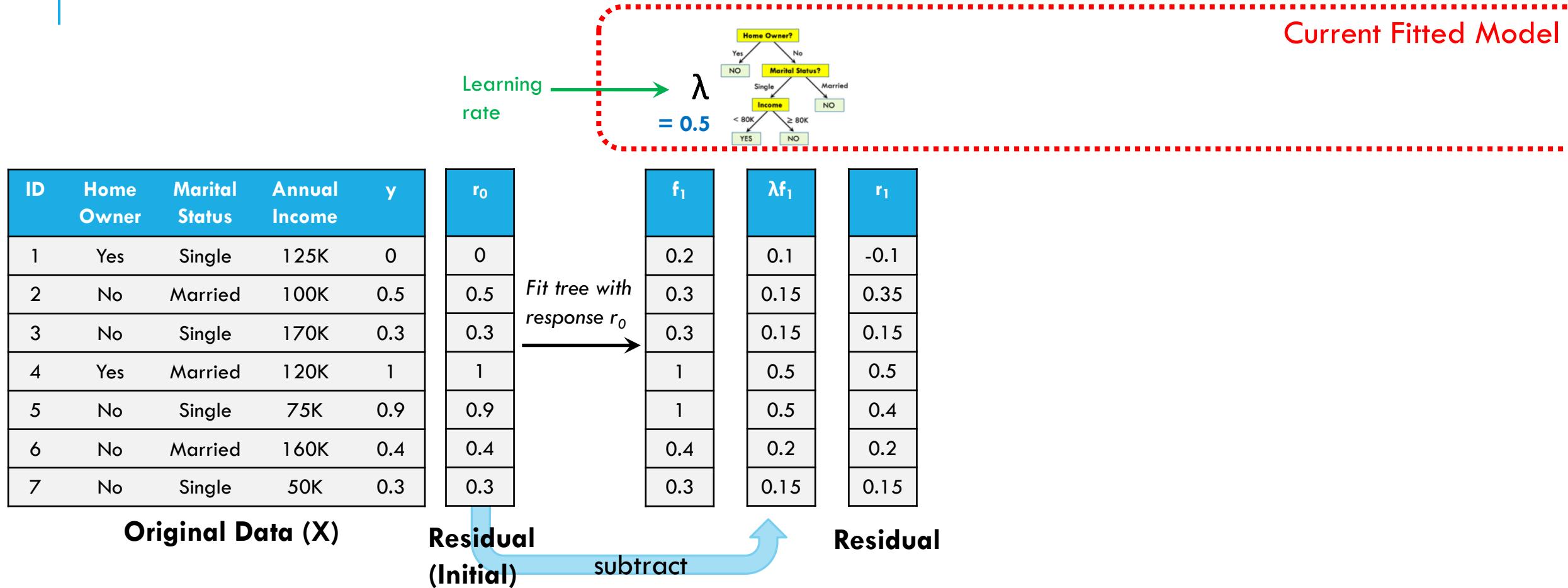


Current Fitted Model



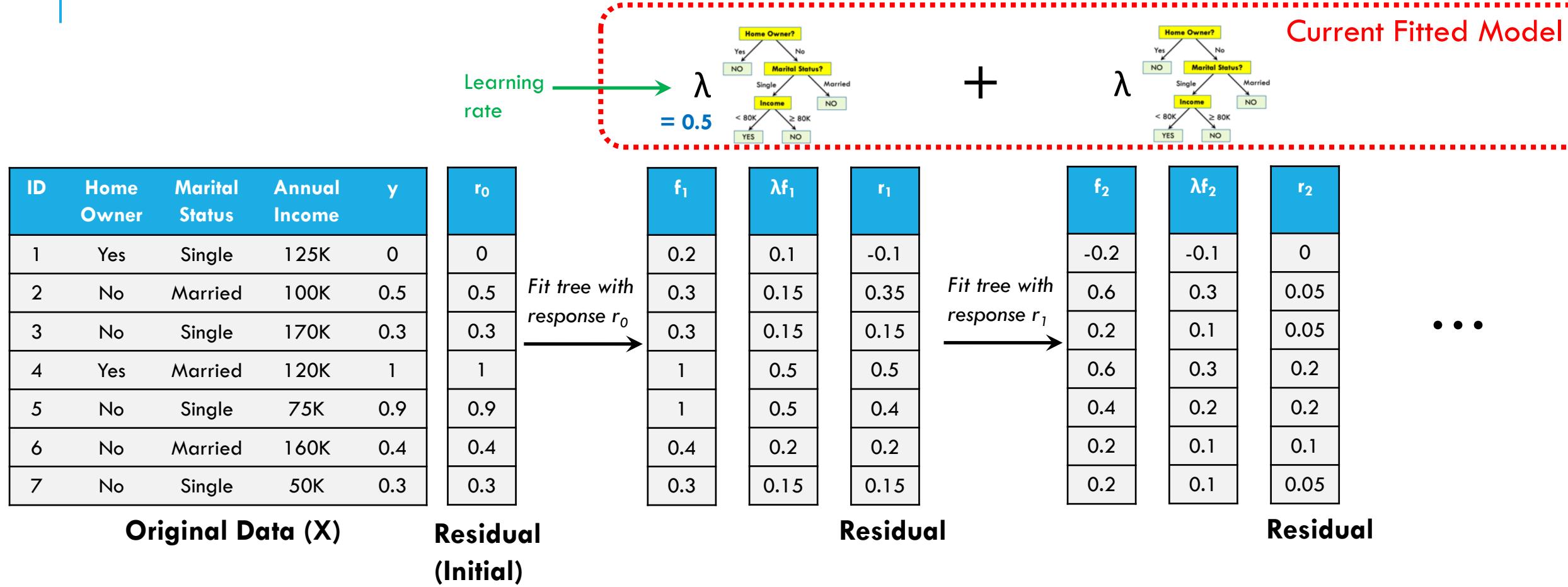
3. Compute 'damped predictions', i.e.  $\lambda$  times decision tree output

# GRADIENT BOOSTED DECISION TREES



4. Compute new residual:  $r_1 = r_0 - \lambda f_1$

# GRADIENT BOOSTED DECISION TREES



5. Repeat the same process (fitting decision tree, etc.) on the new residuals  $r_1$

# GRADIENT BOOSTED DECISION TREES

ID	Home Owner	Marital Status	Annual Income	y
1	Yes	Single	125K	0
2	No	Married	100K	0.5
3	No	Single	170K	0.3
4	Yes	Married	120K	1
5	No	Single	75K	0.9
6	No	Married	160K	0.4
7	No	Single	50K	0.3

Original Data (X)  
Residual  
(Initial)

r <sub>0</sub>
0
0.5
0.3
1
0.9
0.4
0.3

Fit tree with response r<sub>0</sub> →

f <sub>1</sub>
0.2
0.3
0.3
1
1
0.4
0.3

$$\lambda = 0.5$$

$\lambda f_1$
0.1
0.15
0.15
0.5
0.5
0.2
0.15

r <sub>1</sub>
-0.1
0.35
0.15
0.5
0.5
0.2
0.15

Fit tree with response r<sub>1</sub> →

f <sub>2</sub>
-0.2
0.6
0.2
0.6
0.4
0.2
0.2

$$\lambda$$

$\lambda f_2$
-0.1
0.3
0.1
0.3
0.2
0.1
0.1

$$\lambda$$

r <sub>2</sub>
0
0.05
0.05
0.05
0.2
0.2
0.1
0.05

Residual

6. After fitting B trees, the final prediction is the sum of damped trees

Current Fitted Model

+ ...

• • •

# SUMMARY: GRADIENT TREES BOOSTING (REGRESSION CASE)

Model  
Predictions

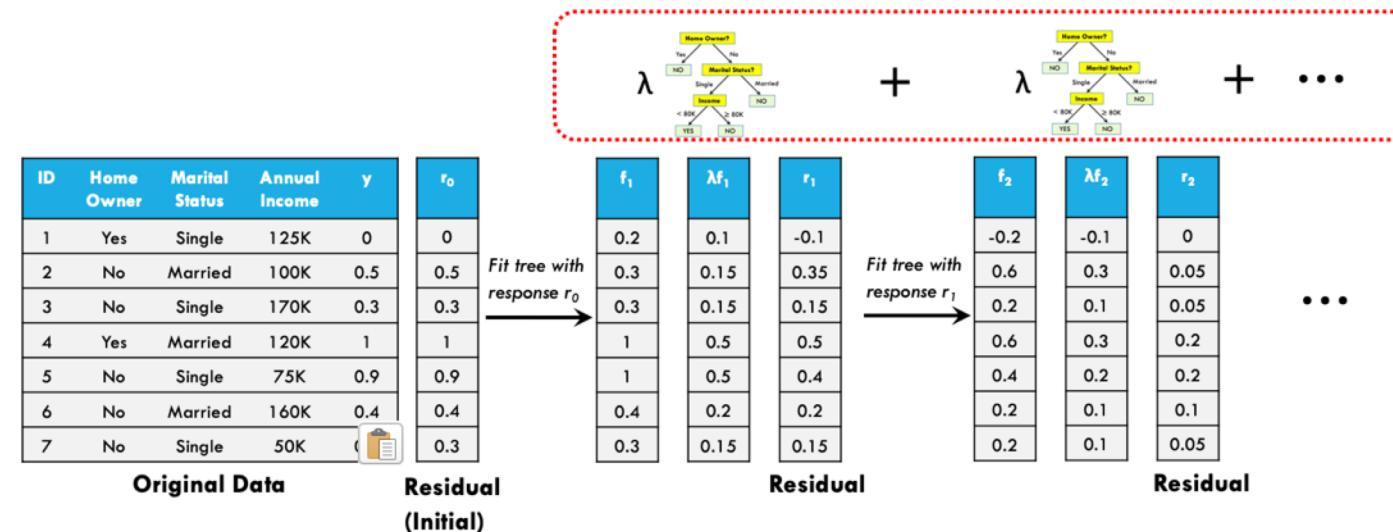
**Initialization:**  $\hat{f}(x) = 0$ , and  $r_i = y_i$  for all  $i$

**For**  $b = 1, \dots, B$ :

- **Fit Decision Tree**  $\hat{f}^b$  to residuals  $r_1, \dots, r_n$ .
- **Update Residual**

$$r_i \leftarrow r_i - \lambda \hat{f}^b(x_i).$$

**Output:**

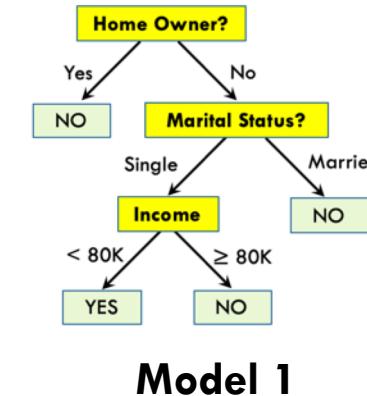
$$\hat{f}(x) = \sum_{b=1}^B \lambda \hat{f}^b(x).$$


# QUIZ: GRADIENT BOOSTING HYPERPARAMETERS



**Q:** While tuning a gradient boosting tree model, you decide to decrease the value of  $\lambda$ . How does this affect the number of trees you should use?

1. Increase
2. Decrease



**Model 1**

• • •



**Model N**

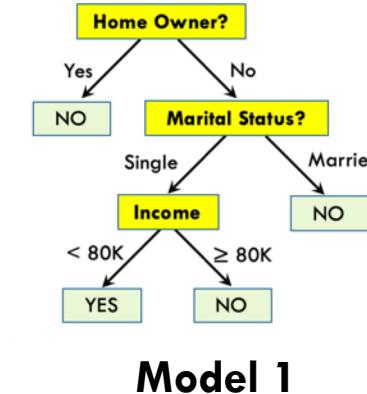
# QUIZ: GRADIENT BOOSTING HYPERPARAMETERS



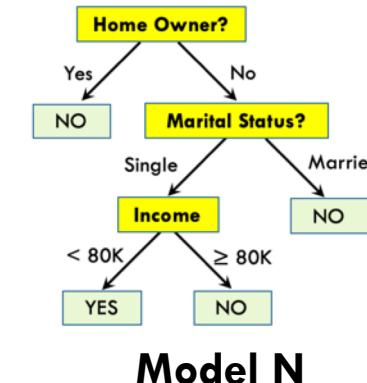
**Q:** While tuning a gradient boosting tree model, you decide to decrease the value of  $\lambda$ . How does this affect the number of trees you should use?

1. Increase 

2. Decrease



• • •



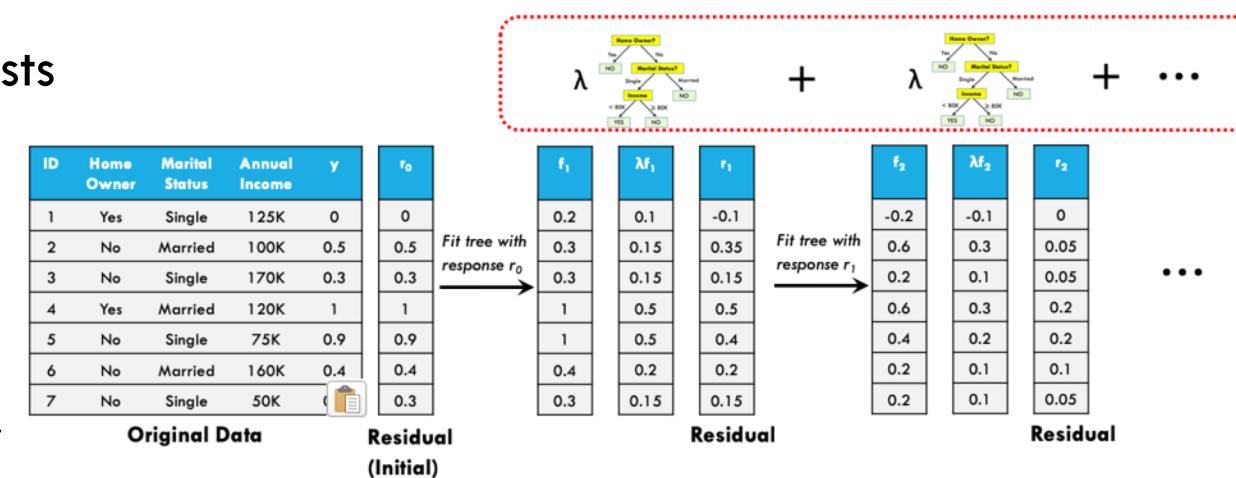
# PROS AND CONS

## Pros

- **Variance Reduction:** Ensembling many decision trees leads to more accurate predictions
- Generally more accurate than random forests

## Cons

- Less interpretable than decision trees
- Slower than decision trees
- Significant tuning required<sup>1</sup> (for tree depth, learning rate, number of trees)



1: **General tuning guidelines:** tree depth is usually between 1 to 10, and is important for controlling overfitting. Tuning it by selecting from {2, 4, 6, 8, 10} based on cross validation is generally fine. Learning rate and no. of trees have to be tuned together: learning rate is usually around 0.01 to 0.1. Smaller learning rate requires more trees (and thus longer training time). If speed is not an issue, a common approach is to set a low learning rate (e.g. 0.01), and train it while selecting the number of trees using early stopping (most gradient boosting tree packages have early stopping already implemented). If this is too slow, it may be advisable to start with a higher learning rate (e.g. 0.1) and later assess whether lowering the learning rate provides improvement in accuracy. See <https://machinelearningmastery.com/configure-gradient-boosting-algorithm/> (and other guides that exist online) for more details (including other tuning parameters)