



Decision Trees , Random Forests and Extra Trees



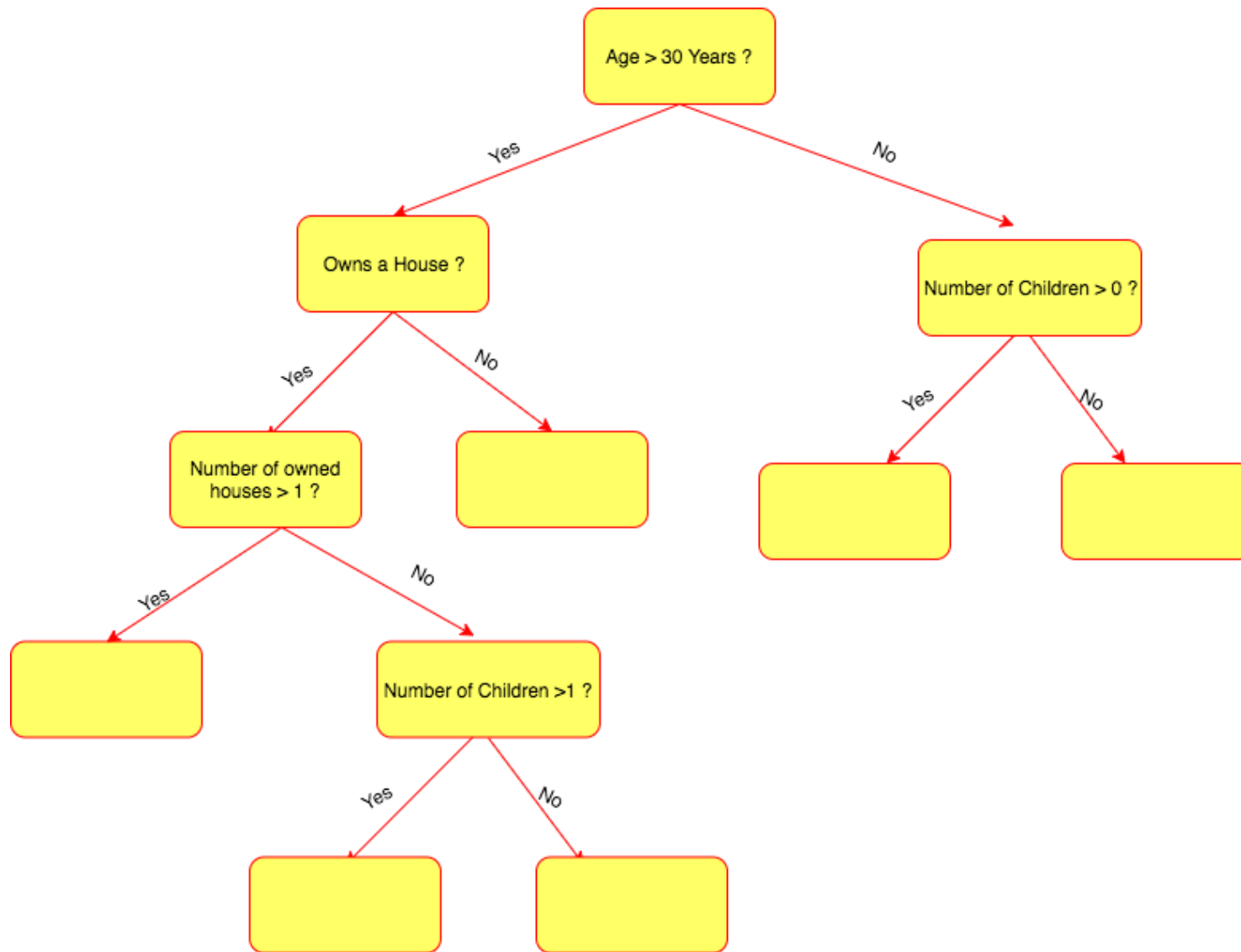
Agenda

Discussion Flow

- Basic Structure of a Decision Tree
- Building a classification/Regression Tree
- Interpretation in absence of parameters/coefficients
- Implementation in Python
- Overfitting issue with Decision Tree
- Random Forests
- Extra Trees
- Implementation in Python

Decision Trees

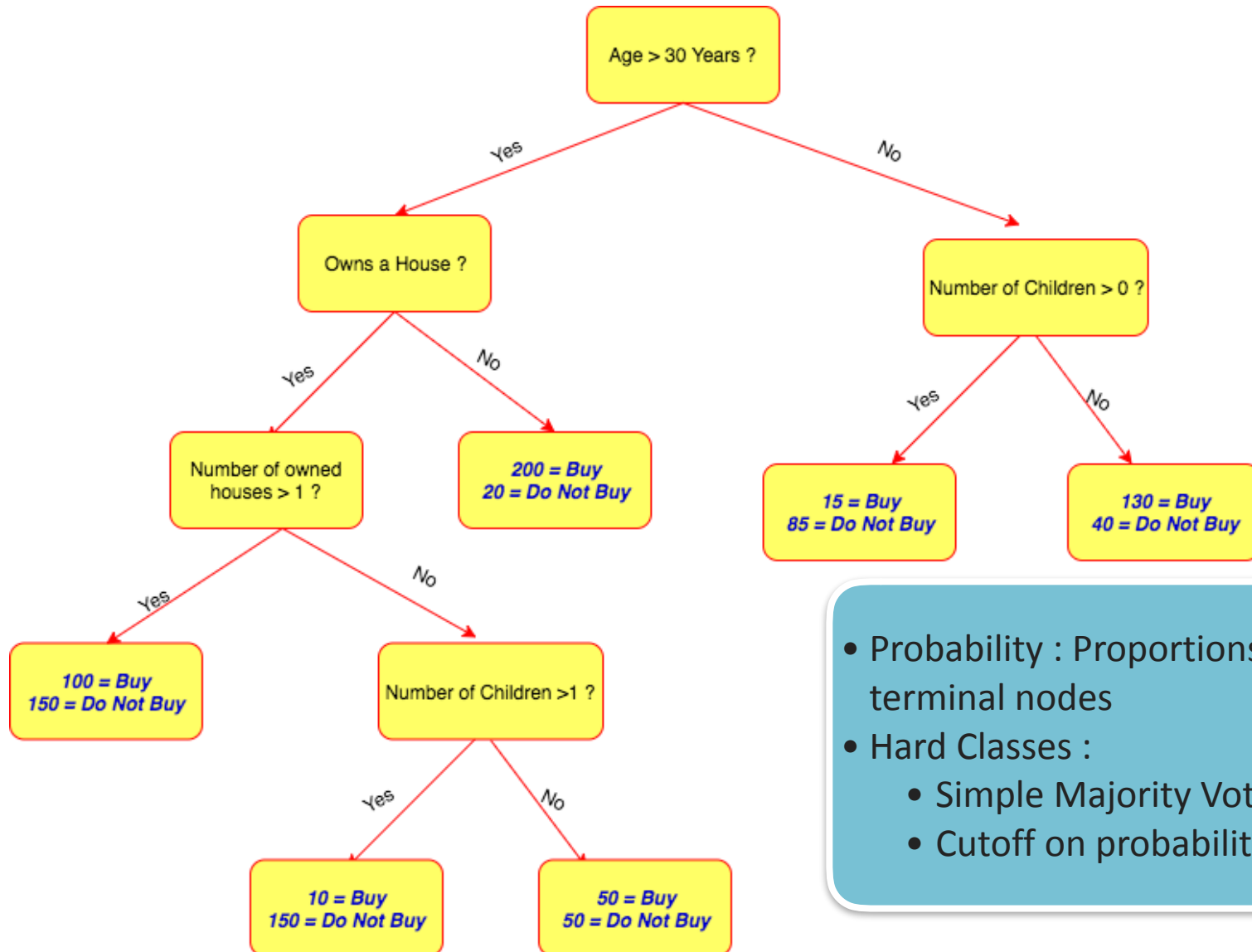
What does it look like?



Questions ? !!

- How do we take decisions ? (Classification/ Regression)
- Where do these rules come from ?
- How do we pick rules for splitting at each node?
- When do we stop splitting nodes ?

How do we take decisions ? (Classification)



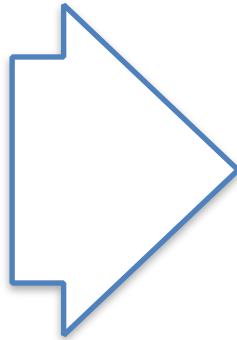
- Probability : Proportions in the terminal nodes
- Hard Classes :
 - Simple Majority Vote
 - Cutoff on probability score

How do we take decisions ? (Regression)

- Much simpler than classification
- Simple Average of target at the terminal node becomes your predicted value for that terminal node

Where do these rules come from? (Numeric Vars)

| Var |
|-----|
| 10 |
| 11 |
| 15 |
| 5 |
| : |
| : |
| : |
| : |
| 24 |



- Range of the numeric variable is broken into intervals
- Based on the values which are there in the data, breaks dont need to be equidistant

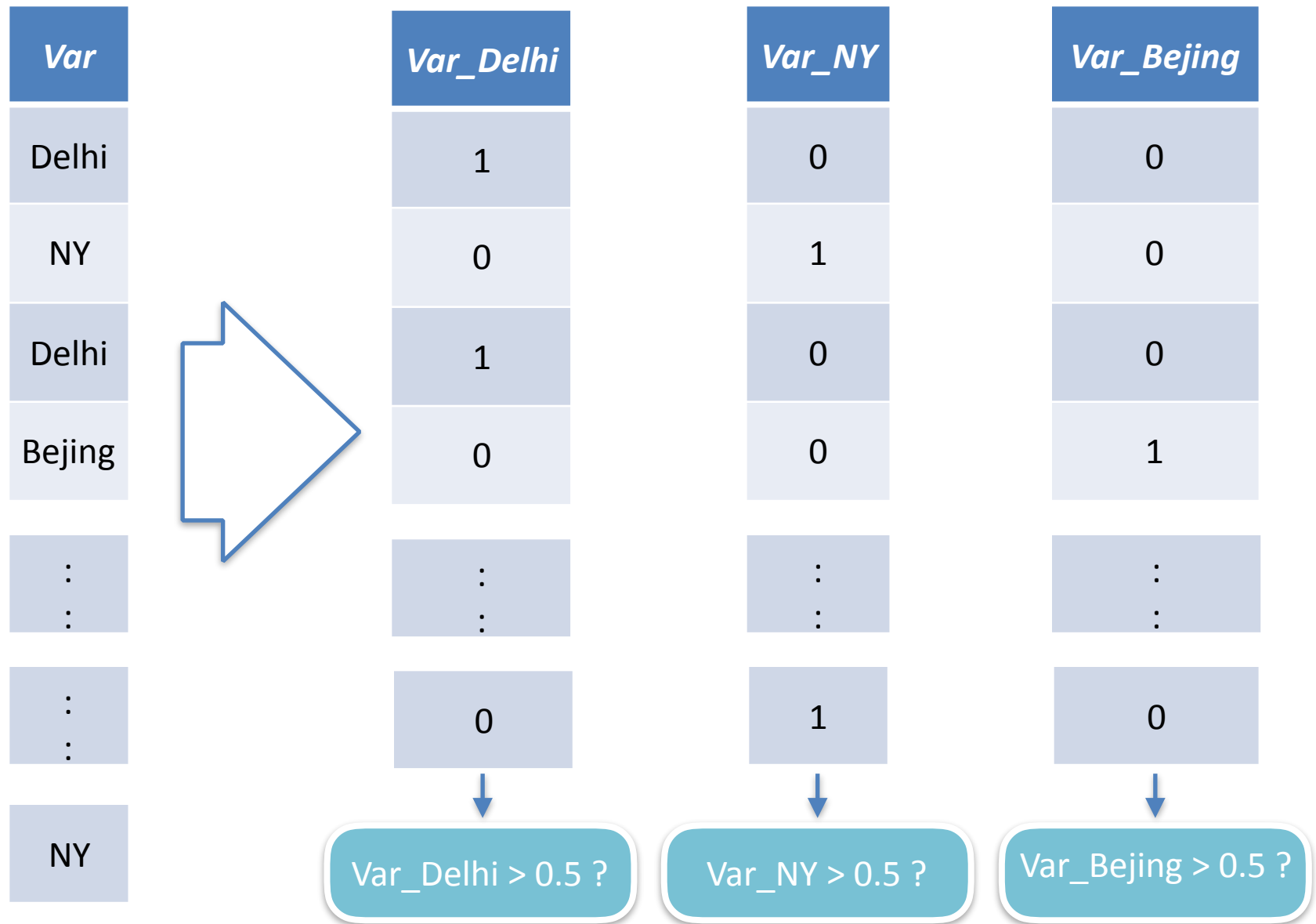
Var > 10 ?

Var > 15 ?

Var > 17.5 ?

Var > 20 ?

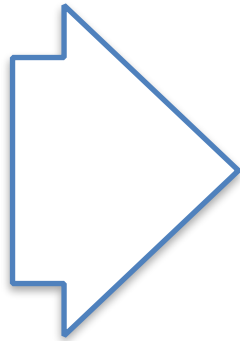
Where do these rules come from? (Categorical vars)



How to Select Rules for split? (Classification)

- Among all the rules available , the one which results in a split with most homogeneous system is selected

Measures of
homogeneity



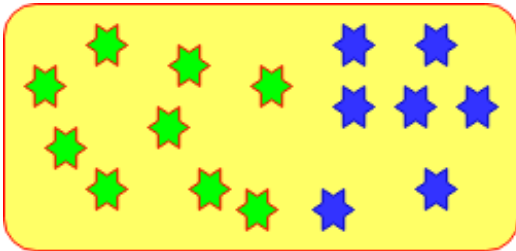
$$gini\ index = 1 - \sum_{i=1}^k p_i^2$$

$$entropy = - \sum_{i=1}^k p_i * \log(p_i)$$

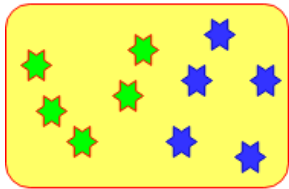
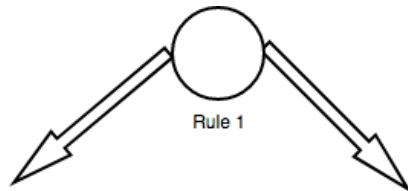
$$deviance = - \sum_{i=1}^k n_i * \log(p_i)$$

Note : There is no theoretical favourite among them , its more of matter of convenience in implementation

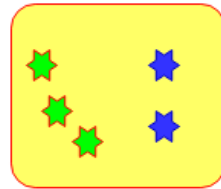
Example: Using entropy for rule selection



$$gini_{parent} = 1 - \left(\left(\frac{8}{15}\right)^2 + \left(\frac{7}{15}\right)^2\right) = .498$$

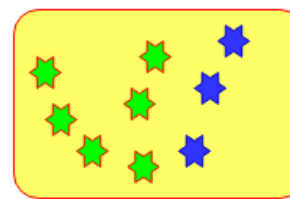
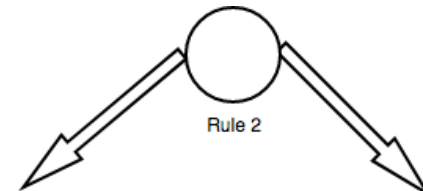


$$gini_{.1} = 1 - \left(\left(\frac{5}{10}\right)^2 + \left(\frac{5}{10}\right)^2\right) = .50$$

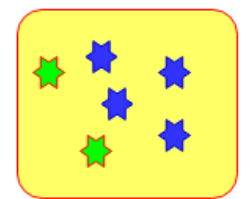


$$gini_{.2} = 1 - \left(\left(\frac{3}{5}\right)^2 + \left(\frac{2}{5}\right)^2\right) = .48$$

$$gini_{new} = \left(\frac{10}{15}\right) * 0.50 + \left(\frac{5}{15}\right) * 0.48 = .493$$



$$gini_{.1} = 1 - \left(\left(\frac{6}{9}\right)^2 + \left(\frac{3}{9}\right)^2\right) = .444$$



$$gini_{.2} = 1 - \left(\left(\frac{2}{6}\right)^2 + \left(\frac{4}{6}\right)^2\right) = .444$$

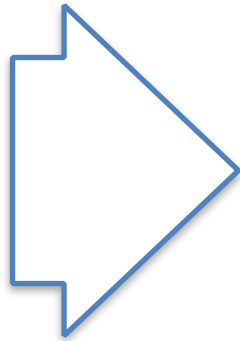
$$gini_{new} = \left(\frac{9}{15}\right) * 0.444 + \left(\frac{6}{15}\right) * 0.444 = .444$$

Rule 2 gets selected for higher decrease in gini

How to Select Rules for split and Make Prediction? (Regression)

- Average of the node is the prediction for the node
- Among all the rules available , the one which results in a split with least Sum of Square of Errors is selected

Measures of
homogeneity



$$SSE = \sum_{i=1}^{n_k} (y_i - \bar{y}_k)^2$$

Example: Using SSE for rule selection

| <i>Response</i> |
|-----------------|
| 5 |
| 6 |
| 12 |
| 11 |
| 4 |
| 8 |
| 13 |
| 5 |
| 6 |
| 7 |

| <i>Error</i> |
|--------------|
| -2.7 |
| -1.7 |
| 4.3 |
| 3.3 |
| -3.7 |
| 0.3 |
| 5.3 |
| -2.7 |
| -1.7 |
| -0.7 |

Prediction

7.7

SSE Parent

92.1

Example Contd Rule 1

| <i>Response</i> |
|-----------------|
| 5 |
| 6 |
| 12 |
| 11 |
| 4 |

| <i>Error</i> |
|--------------|
| -2.6 |
| -1.6 |
| 4.4 |
| 3.4 |
| -3.6 |

| <i>Prediction</i> |
|-------------------|
| 7.6 |

| <i>SSE 1</i> |
|--------------|
| 53.2 |

| <i>SSE new</i> |
|----------------|
| 92 |

| <i>Response</i> |
|-----------------|
| 8 |
| 13 |
| 5 |
| 6 |
| 7 |

| <i>Error</i> |
|--------------|
| 0.2 |
| 5.2 |
| -2.8 |
| -1.8 |
| -0.8 |

| <i>Prediction</i> |
|-------------------|
| 7.8 |

| <i>SSE 2</i> |
|--------------|
| 38.8 |

Example Contd Rule 2

| <i>Response</i> |
|-----------------|
| 5 |
| 6 |
| 4 |
| 5 |
| 6 |

| <i>Error</i> |
|--------------|
| -0.2 |
| 0.8 |
| -1.2 |
| -0.2 |
| 0.8 |

| <i>Prediction</i> |
|-------------------|
| 5.2 |

| <i>SSE 1</i> |
|--------------|
| 2.8 |

| <i>SSE new</i> |
|----------------|
| 29.6 |

| <i>Response</i> |
|-----------------|
| 12 |
| 11 |
| 8 |
| 13 |
| 7 |

| <i>Error</i> |
|--------------|
| 1.8 |
| 0.8 |
| -2.2 |
| 2.8 |
| -3.2 |

| <i>Prediction</i> |
|-------------------|
| 10.2 |

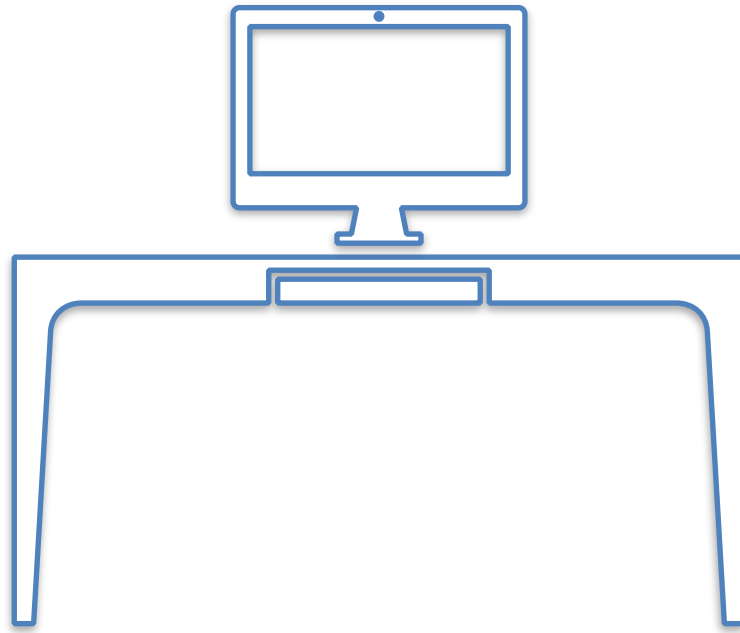
| <i>SSE 2</i> |
|--------------|
| 26.8 |

Rule 2 gets selected
because of higher
decrease in SSE

When do we stop splitting Nodes?

- When does a node become a terminal node?
- When node is completely homogeneous
- When number of observation in the nodes are lower than the specified limit for split
- When all the rules result in a split such that one (or both) child node will have less observation than specified limit for child node
- When number of specified terminal node is reached

Lets see it in action in Python

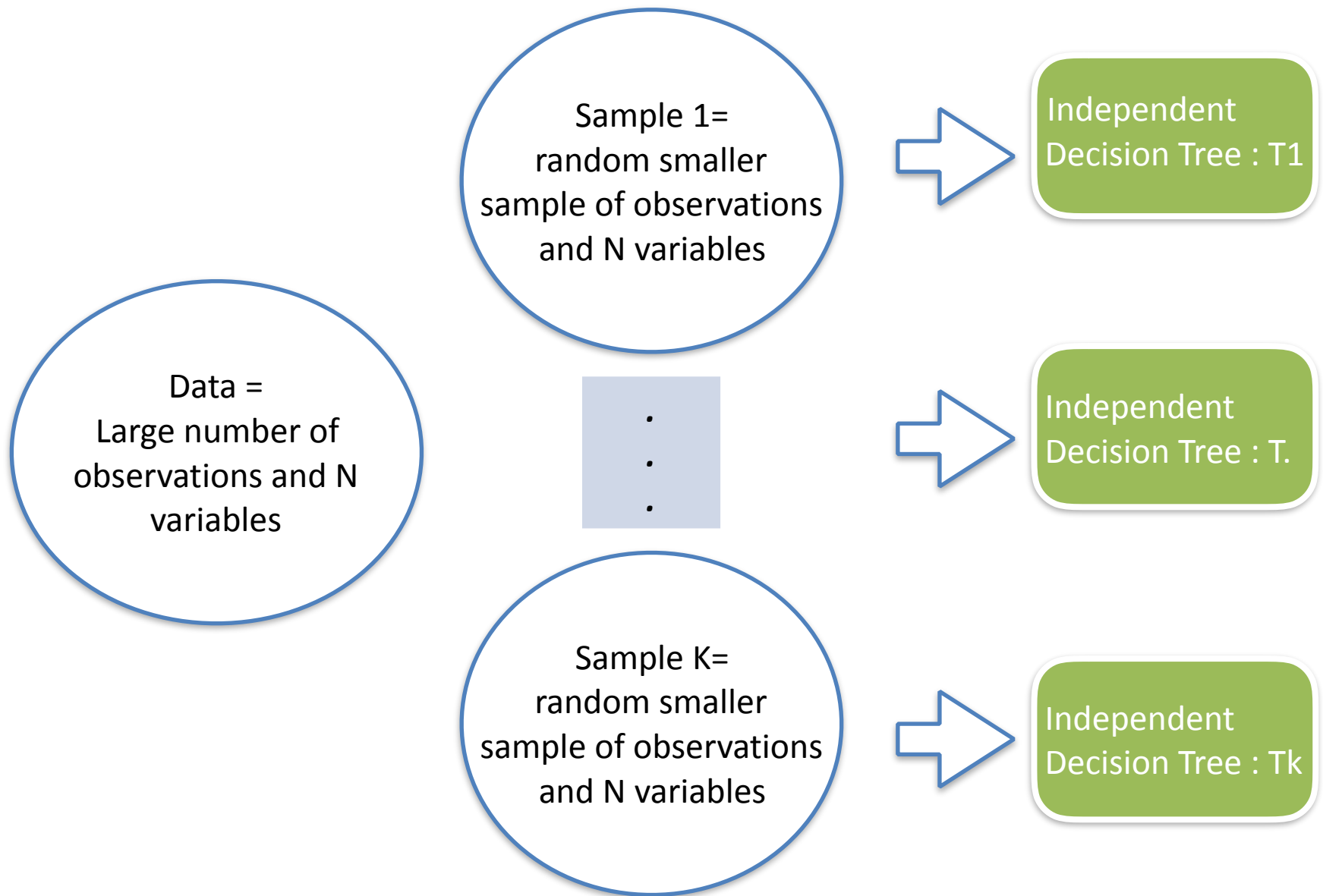


Issues with a single decision tree

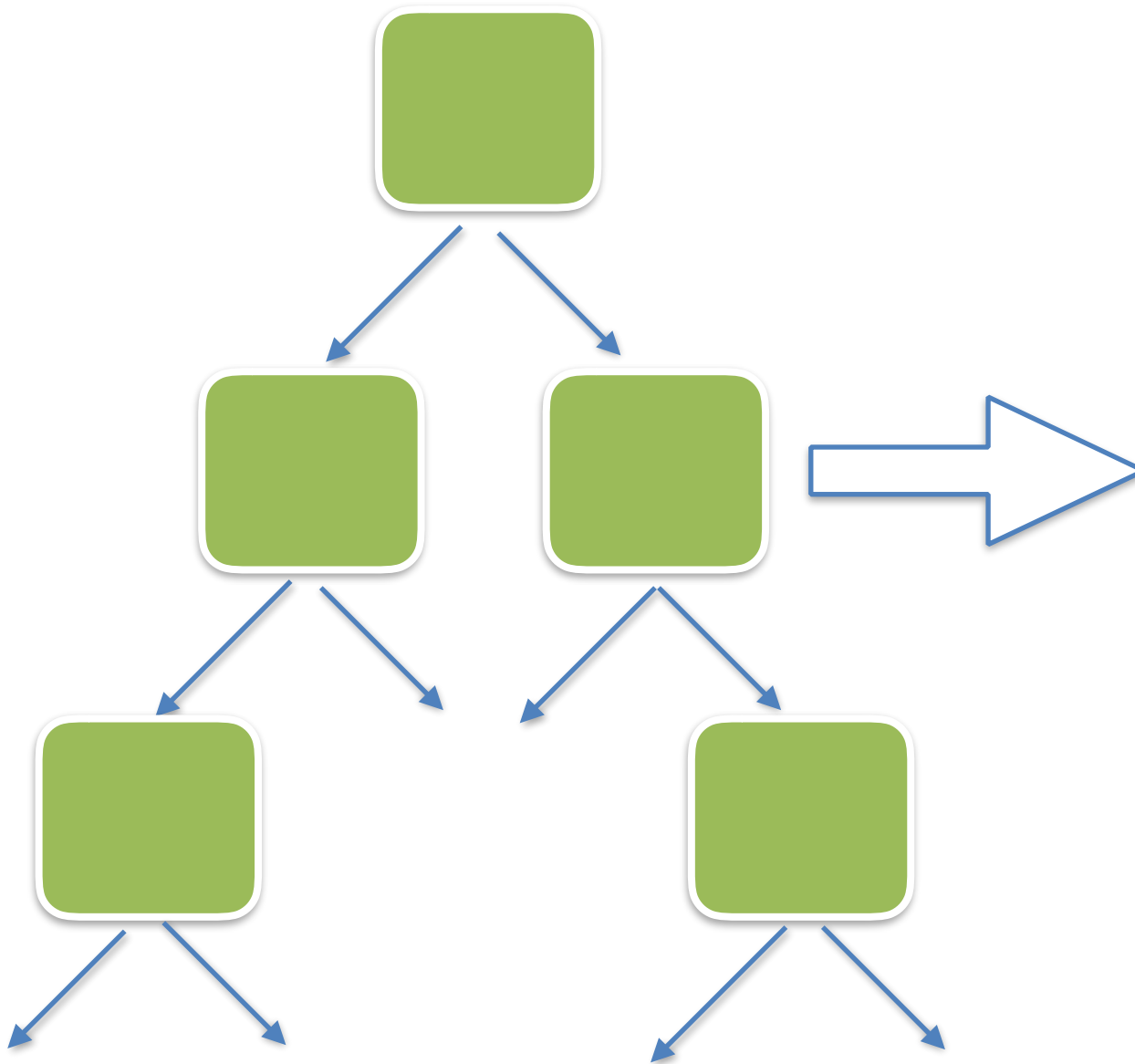
- Susceptible to noisy observations
- Susceptible to noisy variables
- In general overfits the training data

Random Forests

Introduction of randomness in the process



Contd... : for each tree



- At each node(Split) instead of all N variables , only a random subset of variables is considered for splitting
- This random Subset is different for each node across each tree

Extra Trees

Extra Trees

- Short for extremely randomised trees
- Extension of random forest
- In addition to what we do in random forest, extra trees randomly subset rules as well at each node before selecting the best rule for split
- Work well when there are less noisy features but with noisy ranges/categories

Lets see it in action in Python

