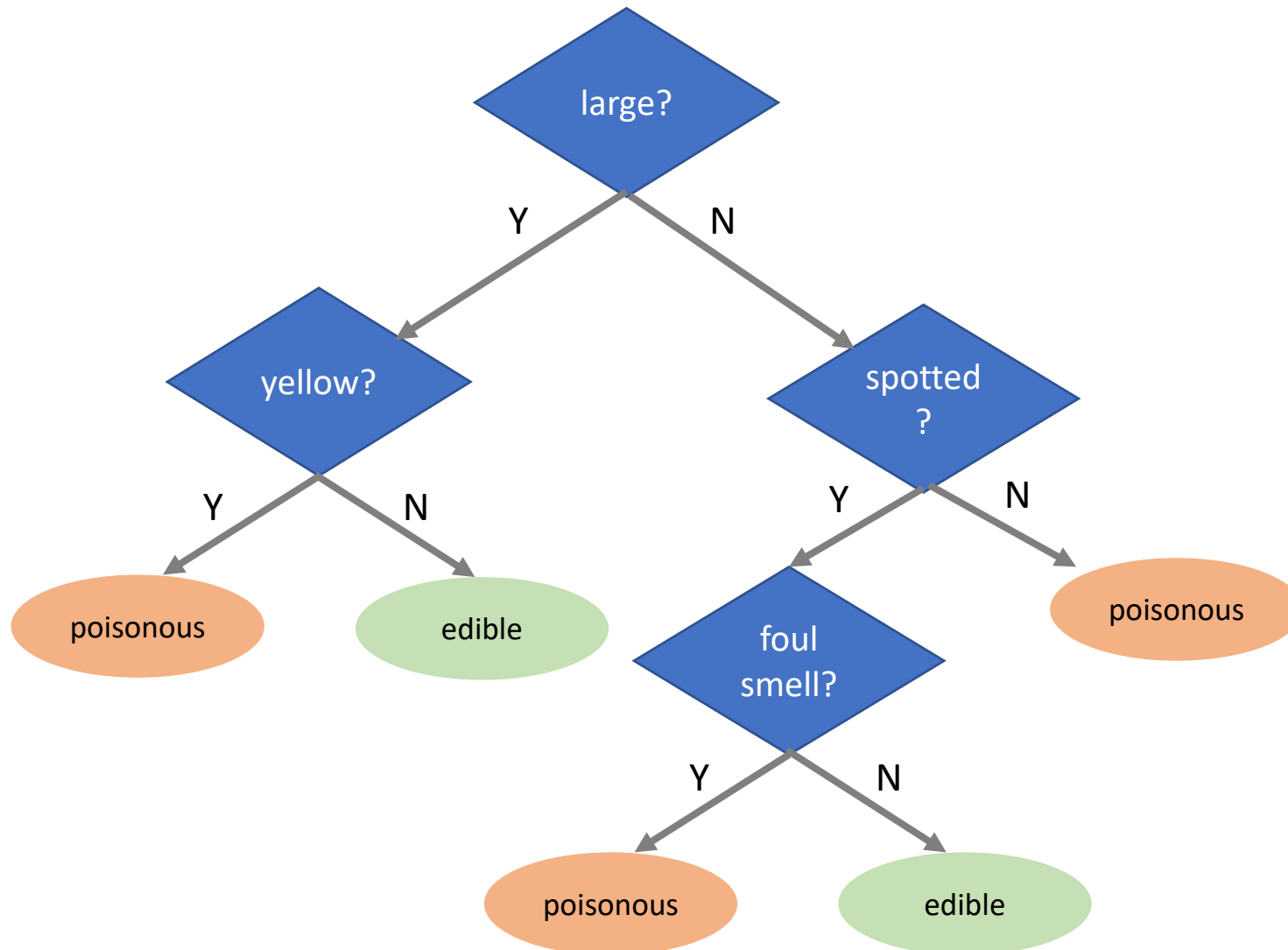# Tree Methods

Geena Kim

# What is Decision Tree?
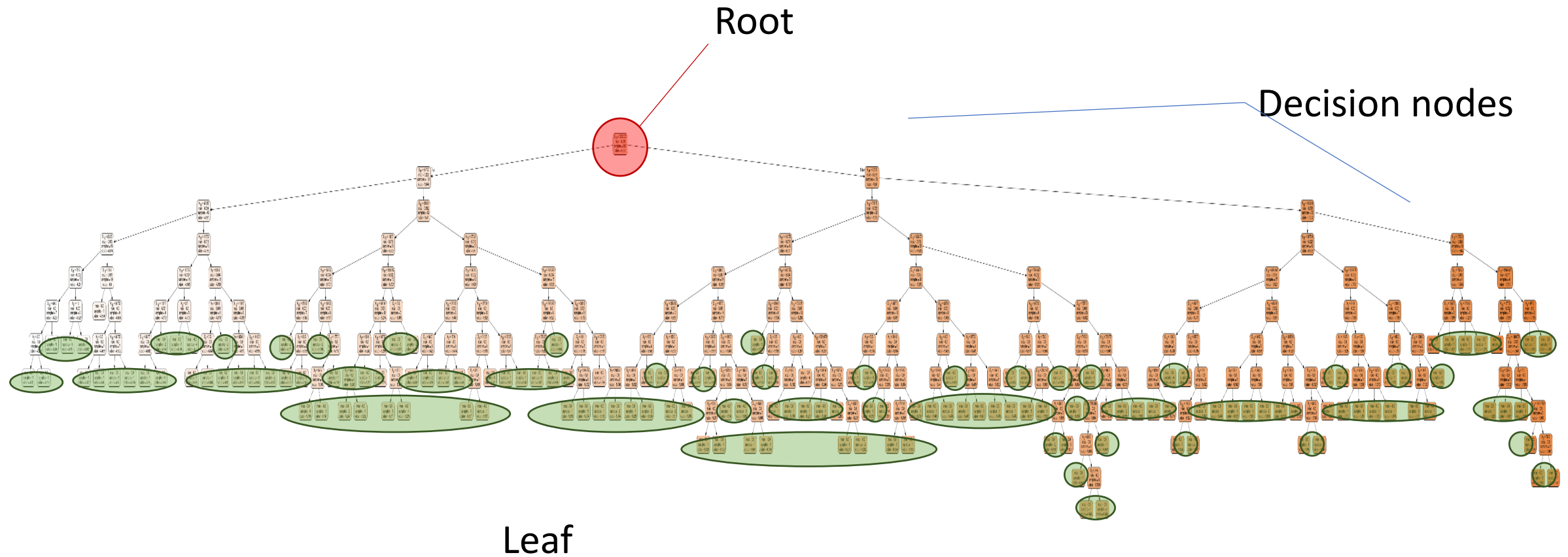


Caesar's mushroom
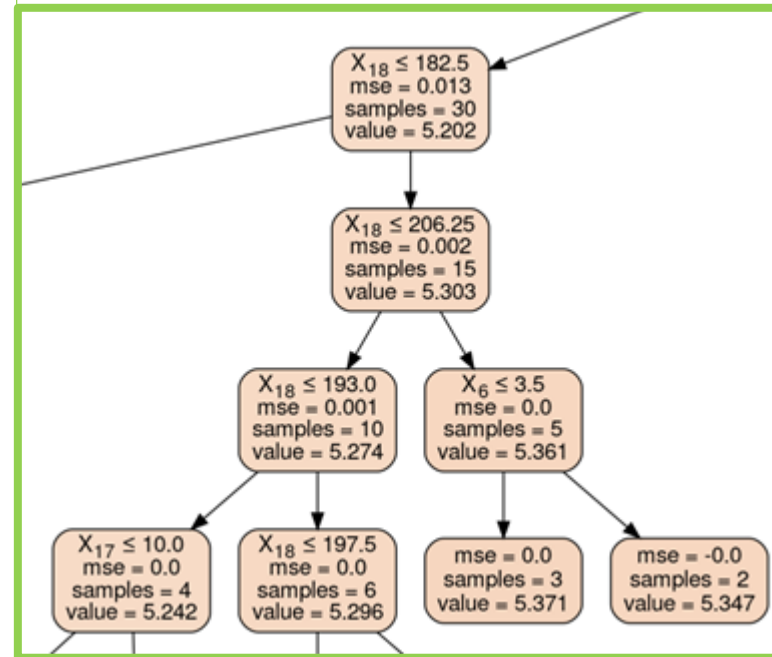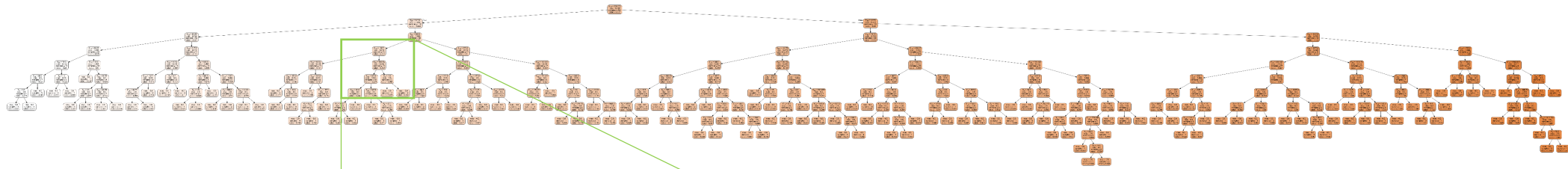


Death Cap

# Decision Tree Nodes



Root

Decision nodes

Leaf

# Decision Nodes



$X_{18} \leq 182.5$
mse = 0.013
samples = 30
value = 5.202

$X_{18} \leq 206.25$
mse = 0.002
samples = 15
value = 5.303

$X_{18} \leq 193.0$
mse = 0.001
samples = 10
value = 5.274

$X_6 \leq 3.5$
mse = 0.0
samples = 5
value = 5.361

$X_{17} \leq 10.0$
mse = 0.0
samples = 4
value = 5.242

$X_{18} \leq 197.5$
mse = 0.0
samples = 6
value = 5.296

mse = 0.0
samples = 3
value = 5.371

mse = -0.0
samples = 2
value = 5.347

# Different kinds of models

Parametric vs. Non-parametric

Parameters vs. Hyperparameters

Linear Regression

Logistic Regression

kNN

Decision Tree

# Optimization objective function

Linear Regression                 Minimize MSE

Logistic Regression           Minimize Cross Entropy

kNN                           No optimization, but uses distance
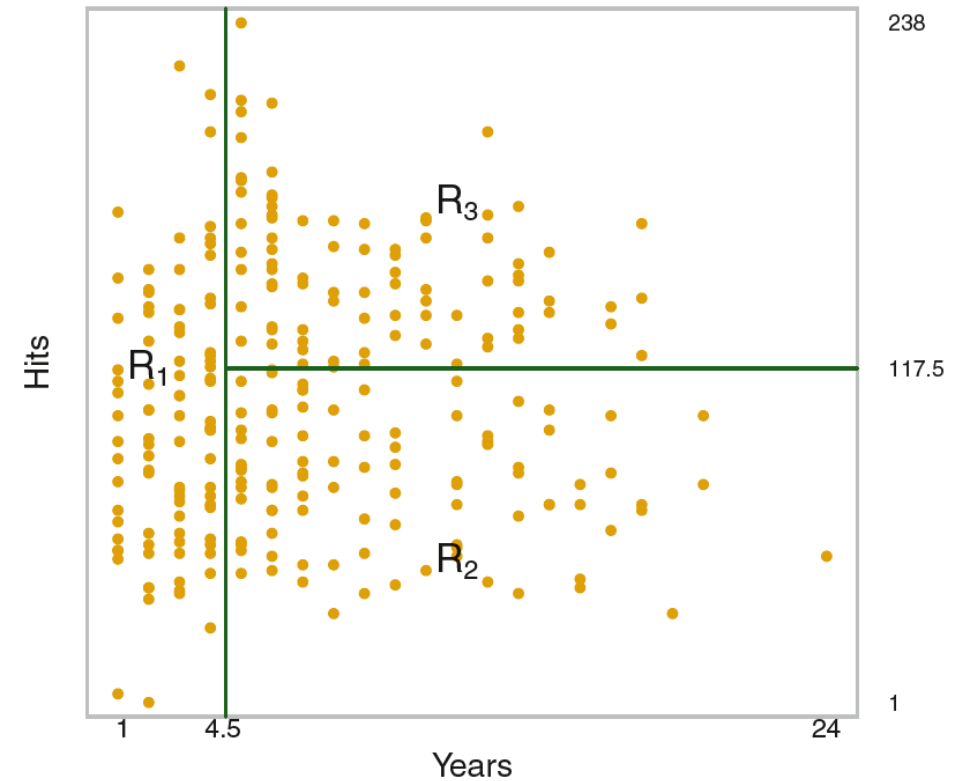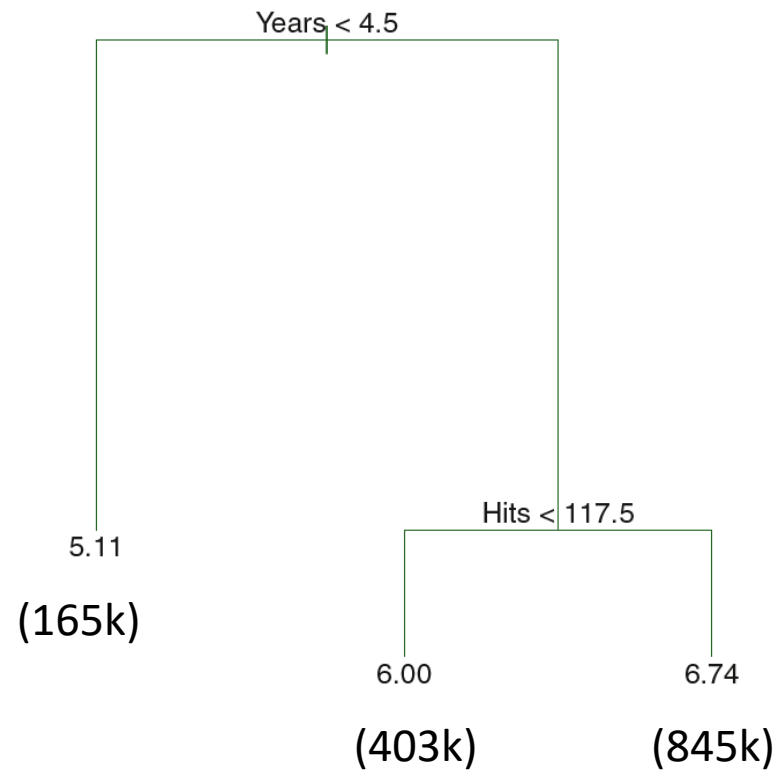
Decision Tree               Split to minimize MSE for Regression task and minimize
Cross Entropy or Gini for Classification task

# Decision Tree Regressor

Predicting Salary of Baseball players
- X1: number of years played in the major league
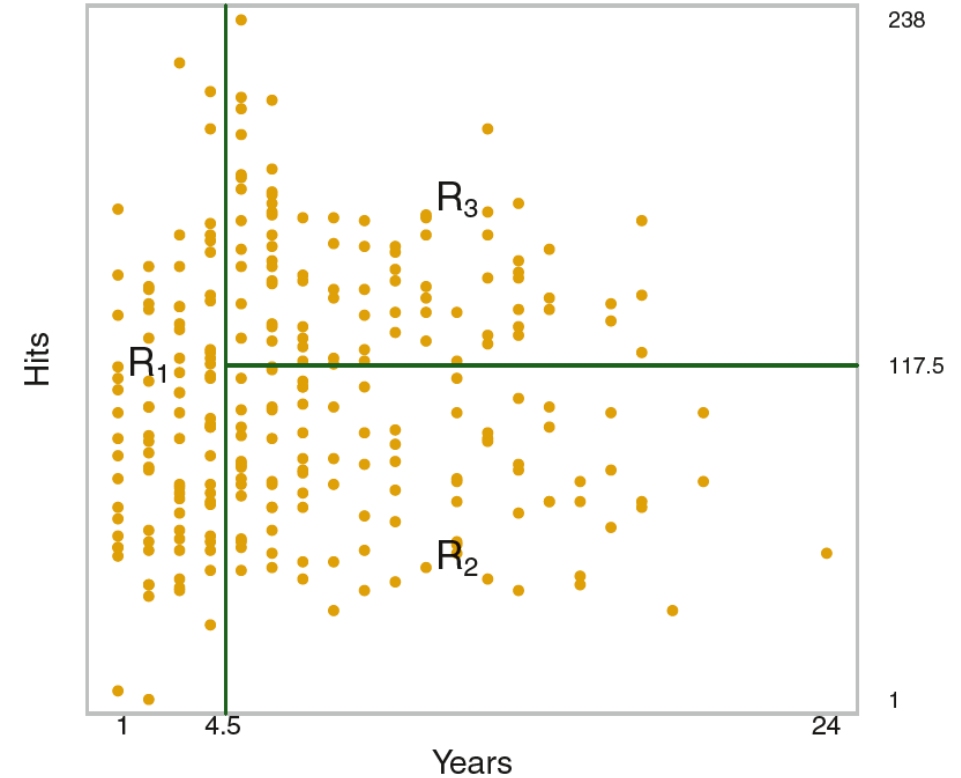- X2: number of hits made in the last year
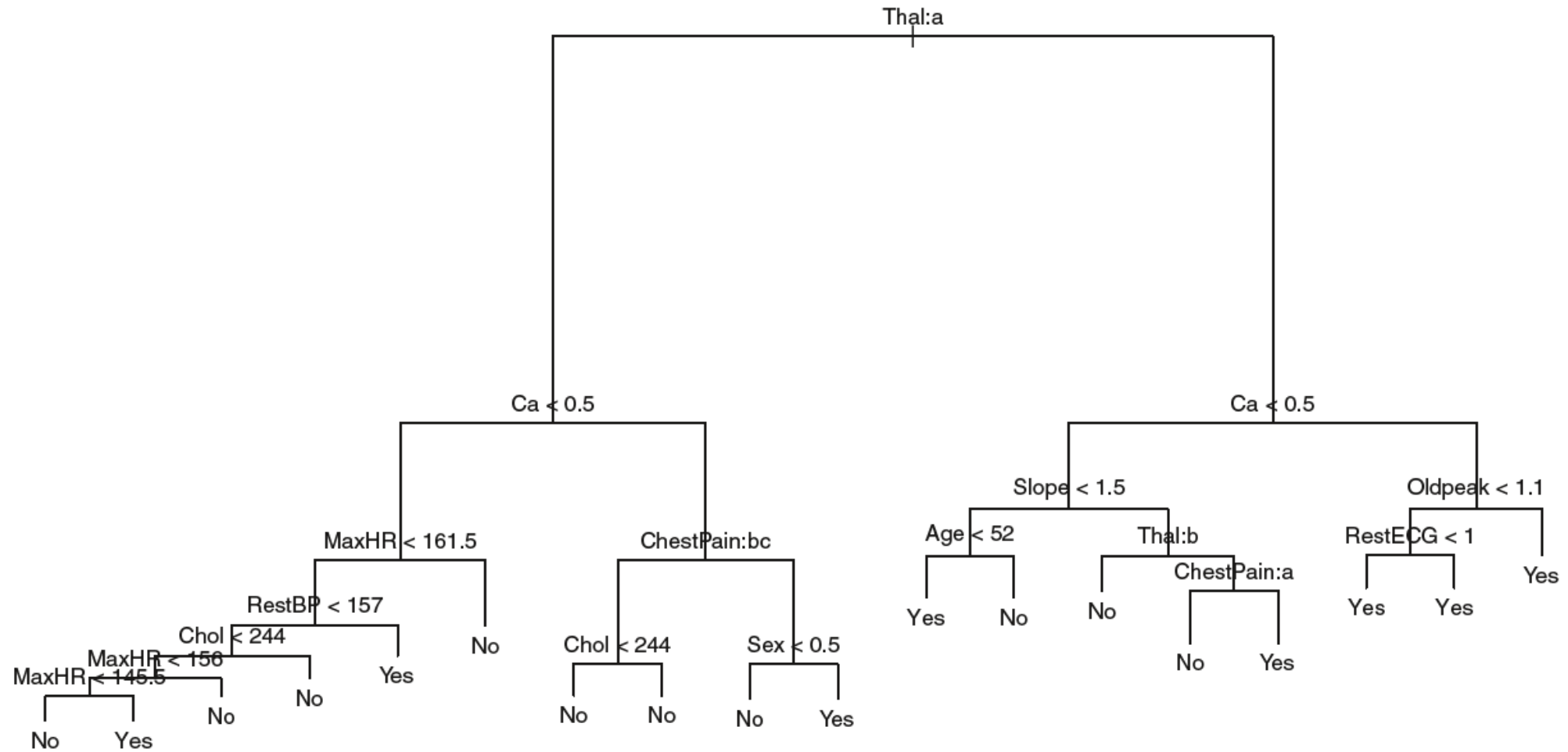- y: log(salary)

# Decision Tree Regressor

The goal is to find boxes R1 ~ RJ such that

$$\sum_{j=1}^{J} \sum_{i \in R_j} \left( y_i - \hat{y}_{R_j} \right)^2$$ is minimized.
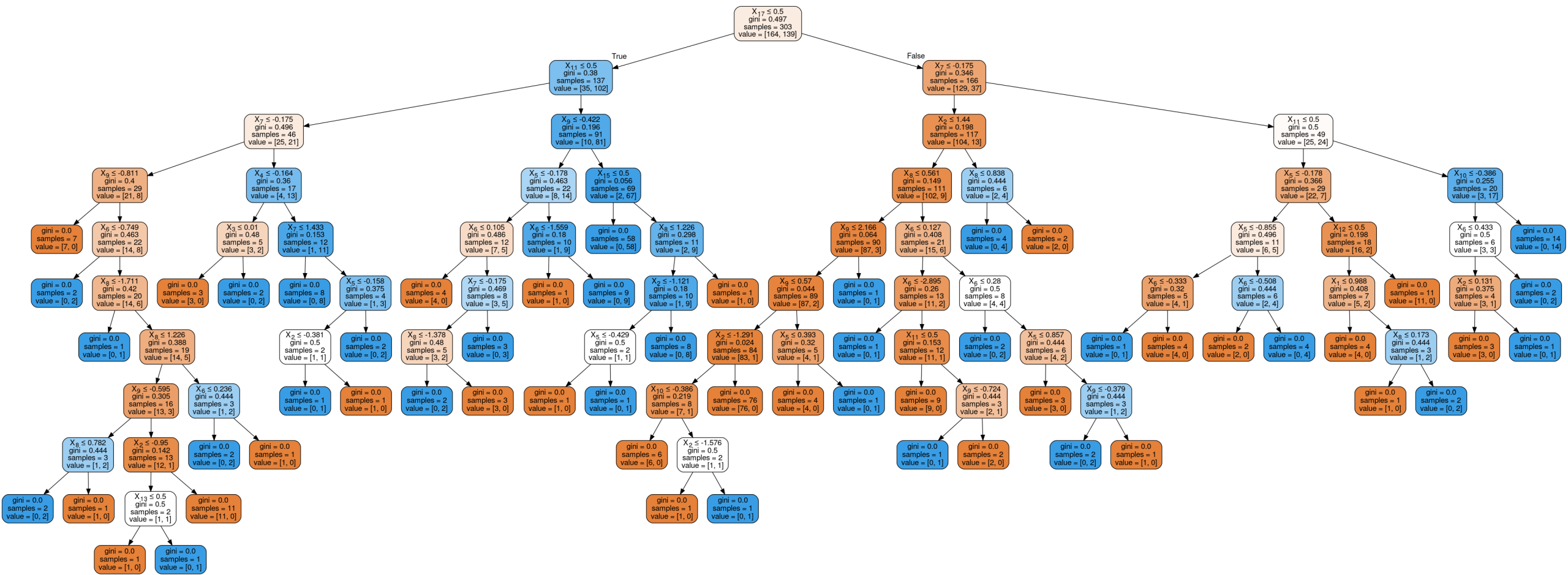
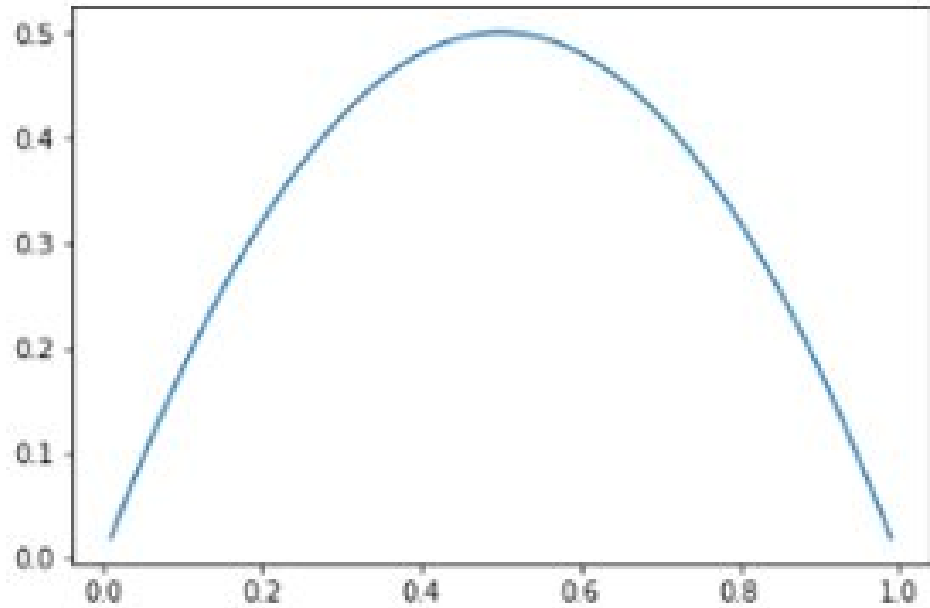the mean of the data in the box

# Decision Tree Classifier

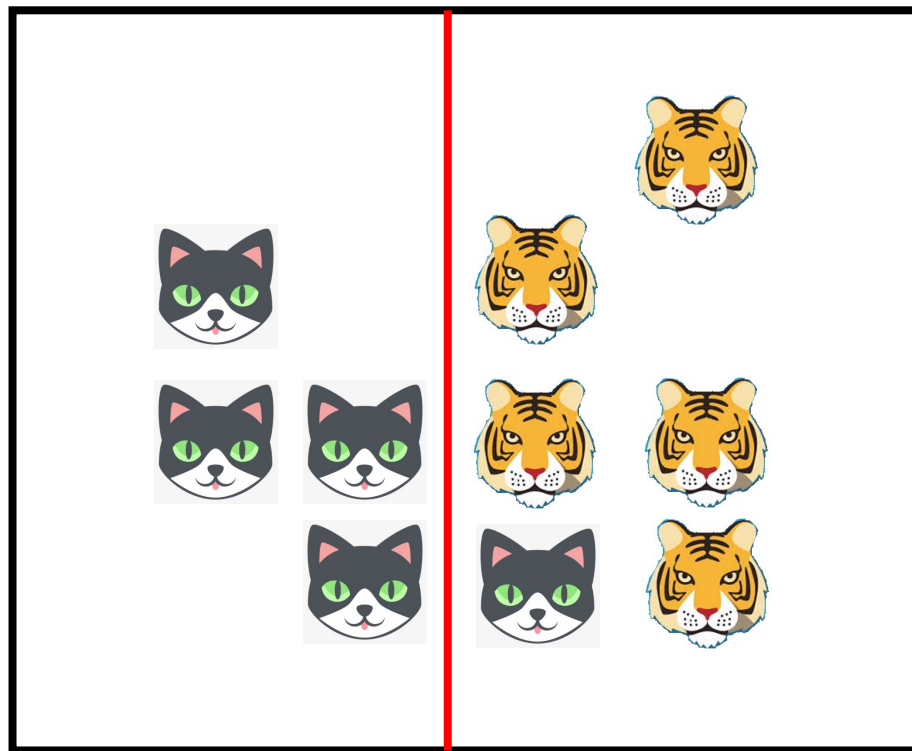# Decision Tree Classifier

# Split criterion- Gini index

```
a = np.arange(0.01,1,0.01)
plt.plot(a,2*a*(1-a));
```
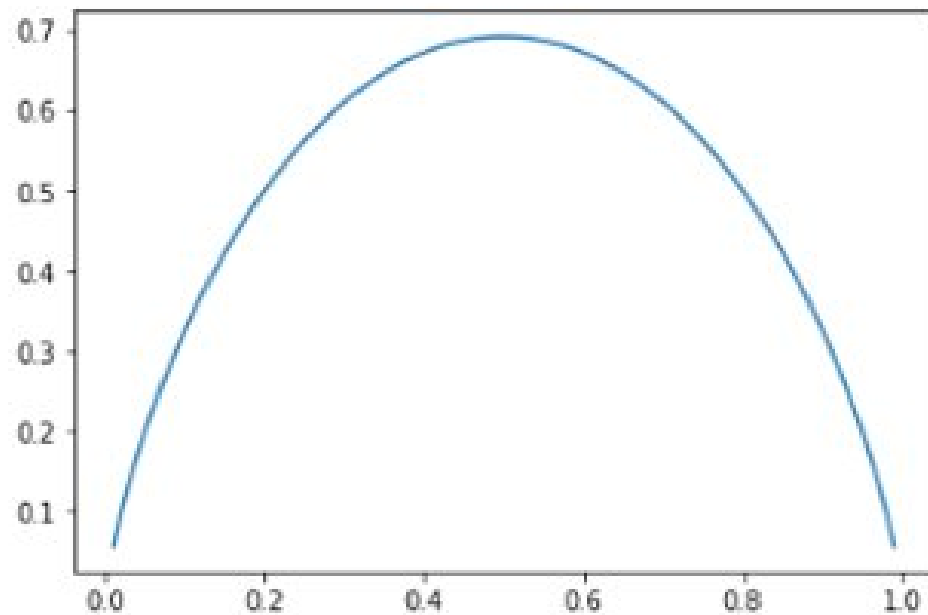


$$H(X_m) = \sum_k p_{mk}(1 - p_{mk})$$

# What is the Gini of this box?

Gini: $H(X_m) = \sum_k p_{mk}(1 - p_{mk})$

# Split criterion- Entropy

```
a = np.arange(0.01,1,0.01)
plt.plot(a,-a*np.log(a)-(1-a)*np.log(1-a));
```
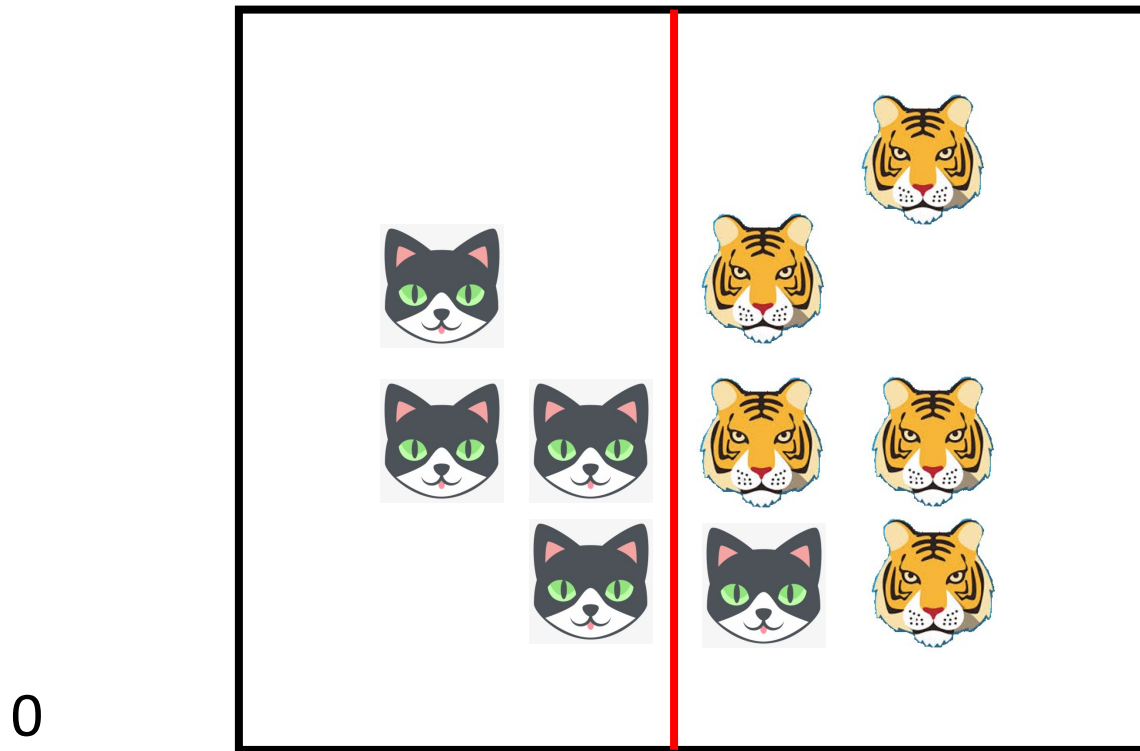


$$H(X_m) = -\sum_k p_{mk} \log(p_{mk})$$

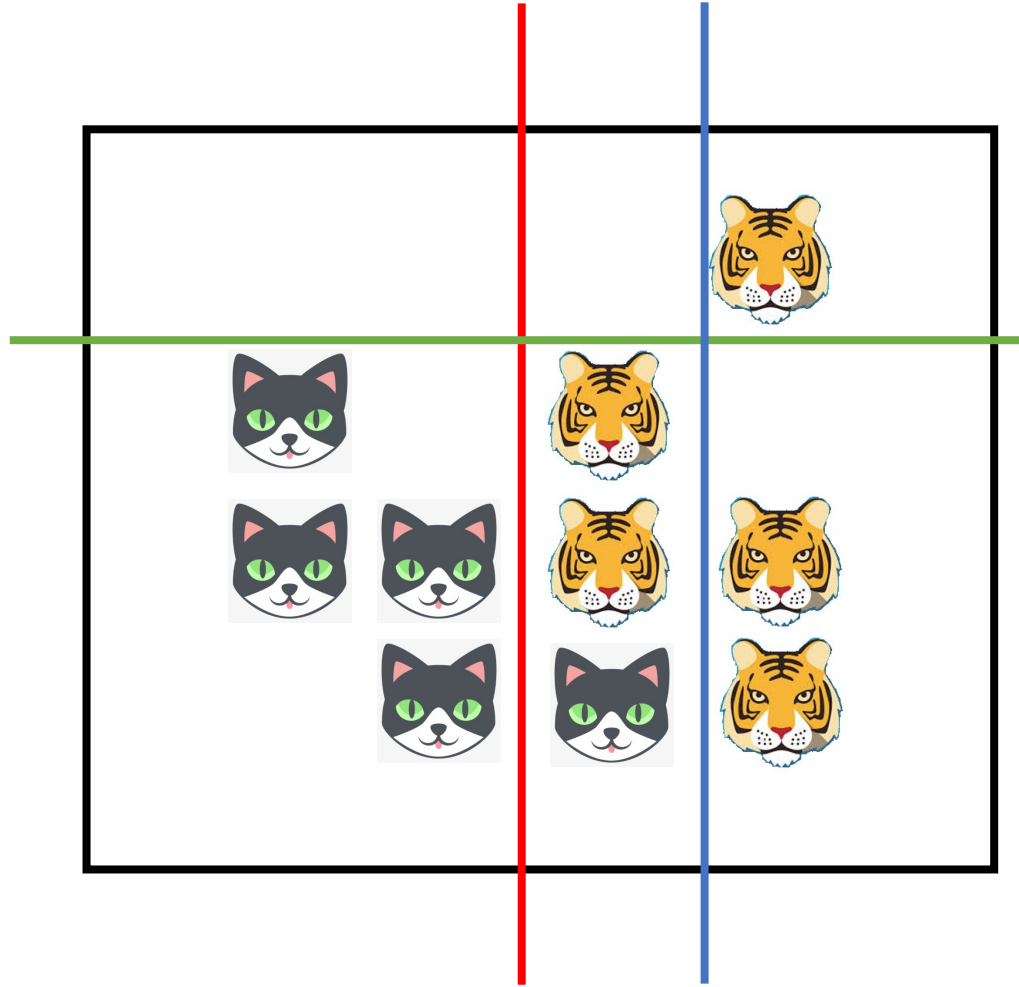# Split criterion- Information gain

Information Gain  = Reduction in Entropy

$$-\left(\frac{1}{2}\log_2\frac{1}{2}+\frac{1}{2}\log_2\frac{1}{2}\right) = 1$$



0

$$-\left(\frac{1}{6}\log_2\frac{1}{6}+\frac{5}{6}\log_2\frac{5}{6}\right) = 0.65$$

Information Gain  = 1 - 0.4*0 - 0.6*0.65 = 0.61

Which split gives the maximum information gain?

# Decision Tree Split Criteria

## Regression Tree

MSE

$$H(X_m) = \frac{1}{N_m} \sum_{i \in N_m} (y_i - \bar{y}_m)^2$$

MAE

$$H(X_m) = \frac{1}{N_m} \sum_{i \in N_m} |y_i - \bar{y}_m|$$

## Classification Tree

Gini

$$H(X_m) = \sum_k p_{mk}(1 - p_{mk})$$

Entropy

$$H(X_m) = -\sum_k p_{mk} \log(p_{mk})$$

Information Gain = E(parent)-E(children)

# Decision Tree – When to stop split?

**max_depth**     The maximum depth of the tree

**min_samples_split**   The minimum number of samples required to split an internal node
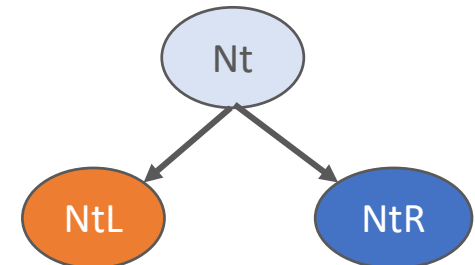
**min_samples_leaf**   The minimum number of samples required to be at a leaf node

**max_features**   The number of features to consider when looking for the best split

**min_impurity_decrease**     A node will be split if this split induces a decrease of the
                                          impurity greater than or equal to this value

The weighted impurity decrease equation is the following:

```
N_t / N * (impurity - N_t_R / N_t * right_impurity
                    - N_t_L / N_t * left_impurity)
```

# Hyperparameter search

## Grid Search Tip

- Give a range of values for each hyperparameter
- Measure a training time for one, then estimate how long for the loop
- Adjust number of values, range, or hyperparameters to include

**max_depth**

**min_samples_split**

**min_samples_leaf**

**max_features**

**min_impurity_decrease**

# Decision Tree Pros and Cons

Trees are easy to understand

Trees don't suffer collinearity

Trees are good for non-linear features

Trees handle categorical variables easily

Trees are weak-learner

Trees have high variance in general

Linear regression is a better choice if features are linear

Tree's performance can be greatly improved when ensembled