# Decision Tree

- can be used for classification & Regression

## Decision Tree Classifier



- Age ← Root Node
- <18, ≥18
- Decision node
- Weight, Smoker
- <60, >60, NO, Yes ← Branch
- Low risk, High risk, Low risk, High risk

Predicted outcomes (Leaf nodes)

## Decision Tree Regressor

mouse eat special diet

T, F

$150 \leq x \leq 180$    $x < 150$

## Classification Trees

- can use both continuous & discrete features
- Numeric threshold can be different for same data

  eg:    T    F
  
  Exercise < 20 min    Exercise < 30

- Prediction nodes can be repeated
- Default: left branch = True, right branch = False

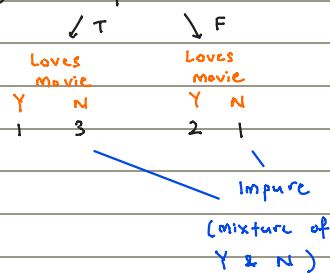## Building classification Trees

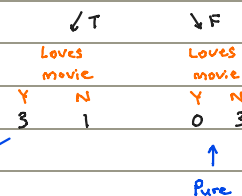1) Determine the feature at root node

  ↳ Steps: Build simple tree for each feature

| Loves Popcorn | Loves Soda | Age | Loves movie |
|---|---|---|---|
| Y | Y | 7 | N |
| Y | N | 12 | N |
| N | Y | 18 | Y |
| N | Y | 35 | Y |
| Y | Y | 38 | Y |
| Y | N | 50 | N |
| N | N | 83 | N |

1) Loves Popcorn feature.          2) Loves Soda feature

  T      F              T      F

Loves Movie    Loves movie      Loves movie    Loves movie

Y   N      Y   N        Y   N      Y   N
1   3      2   1        3   1      0   3

Impure    Pure

(mixture of Y & N)

## Quantifying Impurity 4 leaf nodes

- Gini Impurity (popular)
- Entropy
- Information Gain

## Gini Impurity (categorical)

$$G = 1 - (\text{Probability Yes})^2 - (\text{Probability No})^2$$

$$G\left(\begin{smallmatrix}\text{Loves}\\\text{popcorn}\\(T)\end{smallmatrix}\right) = 1 - \left(\frac{1}{1+3}\right)^2 - \left(\frac{3}{1+3}\right)^2$$

$$= 0.375$$

$$G\left(\begin{smallmatrix}\text{Loves}\\\text{popcorn}\end{smallmatrix} F\right) = 1 - \left(\frac{2}{1+2}\right)^2 - \left(\frac{1}{1+2}\right)^2$$

$$= 0.444$$
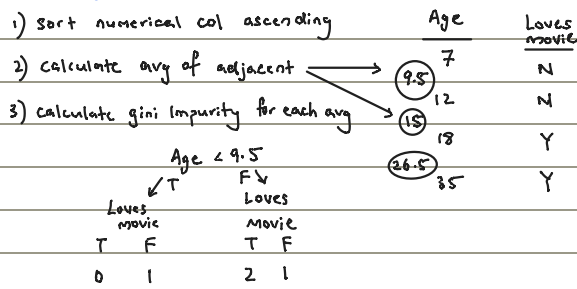
<u>Weighted Gini Impurity</u>

· To account for different sample size

$$G\begin{pmatrix} \text{Loves} \\ \text{popcorn} \end{pmatrix} = \frac{4}{4+3}(0.375) + \frac{3}{4+3}(0.444)$$

$$= 0.405$$

<u>Gini Impurity (numerical)</u>

1) Sort numerical col ascending

2) Calculate avg of adjacent

3) Calculate gini impurity for each avg

| Age | Loves movie |
|---|---|
| 7 | N |
| 12 | N |
| 18 | Y |
| 35 | Y |

(9.5)  (15)  (26.5)

```
        Age < 9.5
       T↙      F↘
   Loves        Loves
   movie        movie
   T   F        T   F
   0   1        2   1
```

$$G\begin{pmatrix} \text{Age} < 9.5 \\ (T) \end{pmatrix} = 1 - (0)^2 - \left(\frac{1}{1}\right)^2$$

$$= 0 \quad \uparrow \text{No Impurity}$$

$$G(\text{Age} < 9.5 \ (F)) = 1 - \left(\frac{2}{3}\right)^2 - \left(\frac{1}{3}\right)^2$$

$$= \frac{4}{9}$$

Weighted Avg $= \frac{1}{4}(0) + \frac{3}{4}\left(\frac{4}{9}\right)$
(Age < 9.5)   $= \frac{1}{3}$

* Repeat for all other avg ages.

eg: | Age < 9.5 : 0.33 | ← Pick lowest Impurity

Age < 15 : 0.343     (if > 2 Ages have same

Age < 26.5 : 0.476     Impurities, can choose either)

$G(\text{Love Popcorn}) = 0.405$

$G(\text{Love soda}) = 0.214$ ← choose this as root node

$G(\text{Age} < 9.5) = 0.33$

2) Perform similar splitting & calculate Gini index

3) Assign output values for each leaf node (threshold)
   ↳ Typically whichever category that has most votes.

<u>Handling overfit of Trees</u>

1) Pruning

2) Limit tree growth (require minimum No. of data in each leaf node)
   ↳ Test optimal No. using
      cross validation