

# DECISION TREES

## Gini Impurity

Now

for a set 'S'

if  $[y_1, y_2, \dots, y_n]$  are the labels.

and let there be  $K \leq n$  different kinds of labels

so let  $p(j)$  be the probability that if you randomly select a label from 'S' it will belong to class  $j$ .

$$p(j) = \frac{|\{y_i, y_i = j\}|}{|\{y_i, i=1 \text{ to } n\}|}$$

$$p(j) = \frac{\text{No. of labels of class 'j'}}{\text{Total no. of labels.}}$$

Now suppose we do an experiment of picking 2 things with replacement from 'S'

The probability that both will be of class  $j$  is  $p(j)^2$

The probability that both will be of same class

$$\begin{aligned} \text{is } & p(1)^2 + p(2)^2 + \dots + p(s)^2 \\ & = \sum_{i=1}^j p(i)^2 \end{aligned}$$

Thus the probability of 2 diff groups

is 
$$1 - \sum_{i=1}^j p(i)^2$$

This should be minimum as we don't want various outputs

and This is defined to be gini impurity of set S

$$G(S) = 1 - \sum_{i=1}^j p(i)^2$$

This is same as  $G(S) = \sum_{i=1}^j p(i)(1-p(i))$

as 
$$\sum_{i=1}^j p(i)(1-p(i))$$

$$= \sum_{i=1}^j p(i) - p(i)^2$$

$$= \underbrace{\sum_{i=1}^j p(i)}_{\text{probability of picking any class} = 1} - \sum_{i=1}^j p(i)^2$$

$$= 1 - \sum_{i=1}^j p(i)^2$$