

# Week 1 - Decision Tree

## 0.Entropy

**Def** measurement of random variable's uncertainty

**Equ** empirical entropy  $H(X)$ , empirical conditional entropy  $H(Y|X)$

$$\begin{aligned}
 H(Y|X) &= H(X, Y) - H(X) \\
 &= -\sum_{x,y} P(X, Y) \log P(X, Y) + \sum_x P(X) \log P(X) \\
 &= -\sum_{x,y} P(X, Y) \log P(X, Y) + \sum_x (\sum_y P(X, Y)) \log P(X) \\
 &= -\sum_{x,y} P(X, Y) \log P(X, Y) + \sum_{x,y} P(X, Y) \log P(X) \\
 &= -\sum_{x,y} \log \frac{P(X, Y)}{P(X)} \\
 &= -\sum_{x,y} \log P(Y|X) \\
 &= -\sum_x \sum_y P(X) P(Y|X) \log P(Y|X) \\
 &= -\sum_x P(X) \sum_y P(Y|X) \log P(Y|X) \\
 &= \sum_x P(X) H(Y|X = x_i)
 \end{aligned}$$

**Ann** binary logarithm

## 1. Feature Selection

### 1.1 information gain

**Def** mutual information of data set and feature

**Equ** feature  $A$  in  $(a_1, a_2, \dots, a_n)$ , data set  $D$ , information gain  $g$ , empirical entropy  $H$ , class number  $K$

$$\begin{aligned}
 g(D, A) &= H(D) - H(D|A) \\
 &= -\sum_{k=1}^K P(C_k) \log P(C_k) + \sum_{i=1}^n P(A_i) \sum_{k=1}^K P(D_k|A_i) \log P(D_k|A_i) \\
 &= -\sum_{k=1}^K \frac{|C_k|}{|D|} \log \frac{|C_k|}{|D|} + \sum_{i=1}^n \frac{|D_i|}{|D|} \sum_{k=1}^K \frac{|D_{ik}|}{|D_i|} \log \frac{|D_{ik}|}{|D_i|}
 \end{aligned}$$

### 1.2 information gain ratio

**Def** Normalization: tackling the possibility of a overwhelming variable set

$$\begin{aligned} \text{Equ } g_R(D, A) &= \frac{g(D, A)}{H_A(D)} \\ &= \frac{g(D, A)}{-\sum_{i=1}^n \frac{|D_i|}{|D|} \log \frac{|D_i|}{|D|}} \end{aligned}$$

### 1.3 Gini coefficient

**Def** purity of the data set

$$\begin{aligned} \text{Equ } Gini(D) &= \sum_{i \neq j} P_i P_j \\ \because \text{binaryTree} \\ &= \sum_{k=1}^K P_k (1 - P_k) \\ &= 1 - \sum_{k=1}^K (P_k)^2 \\ &= 1 - \sum_{k=1}^K \left( \frac{|C_k|}{|D|} \right)^2 \end{aligned}$$

$$\text{Equ } Gini(D, A) = \frac{|D_1|}{|D|} Gini(D_1) + \frac{|D_2|}{|D|} Gini(D_2)$$

## 2.Generation

---

**Def**

*ID3* ---- information gain

*C4.5* ---- information gain ratio

- 1) calculate each gain , choose the maximum gain.
- 2) divide recursively until (subtree in same class) or (gain  $\leq$  threshold)

*CART* ---- Gini coefficient

- 1) choose the min  $Gini(D, A_i)$  as the optimal segmentation point
- 2) divide recursively until (each feature traversed) or (subtree in same class)

## 3.Pruning

---

**Def** alleviate degree of overfitting

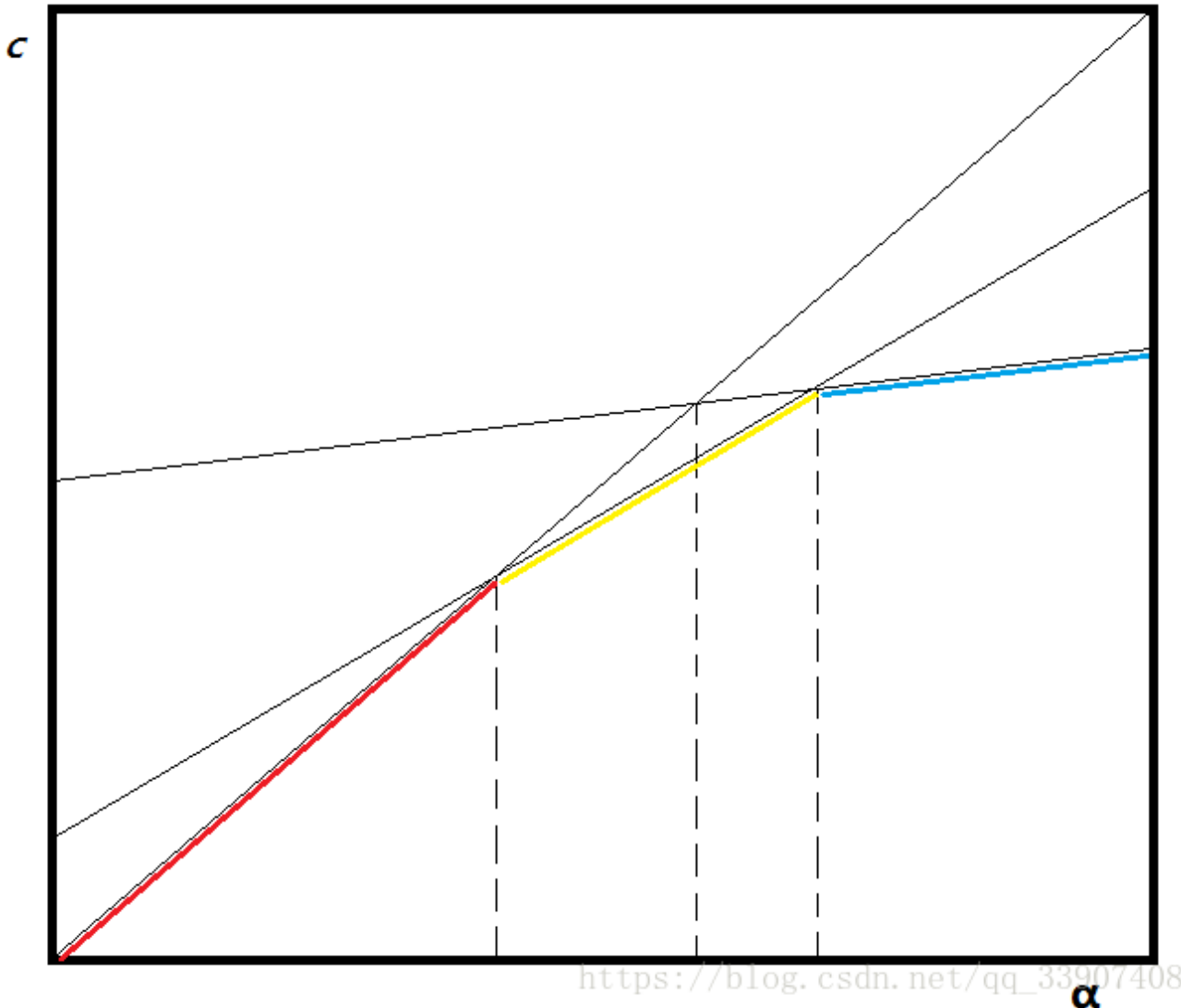
**loss function** a function that maps an event or values of one or more variables onto a real number intuitively representing some “cost” associated with the event.

**Equ** leaf node number  $|T|$ , class number  $N_t$ , parameter  $\alpha$

$$C_\alpha(T) = \sum_{t=1}^{|T|} N_t H_t(T) + \alpha |T|$$

$$= - \sum_{t=1}^{|T|} \sum_{k=1}^K N_{tk} \log \frac{N_{tk}}{N_t} + \alpha |T|$$

decreasing loss ---- need pruning



### *CART Pruning*

Situation: pruning without a given parameter  $\alpha$

In the way of loss function, leafy trees benefit from little  $\alpha$  and the other way around. Pruning happens when loss function of a inner node acting as a root equals to the one acting as a leaf.

$$C_\alpha(t) = C(t) + \alpha$$

$$C_\alpha(T_t) = C(T_t) + \alpha |T|$$

$$\therefore C_\alpha(t) = C_\alpha(T_t)$$

$$\therefore \alpha = \frac{C(t) - C(T_t)}{|T| - 1}$$

- 1) each inner node has an  $\alpha$
- 2) acquire an ascending set of  $\alpha$ , and a set of subtrees accordingly
- 3) Cross Validation: test the subtree set, select the one with the highest accuracy rate

