

Entropy

Now

Consider a set $S = \{y_1, \dots, y_n\}$ n elements

Let K be the no: of different classes

Now let $p_S(i)$ be the probability that if you randomly select a label y' from S it will belong to class ' i '

$$\text{so } p_S(i) = \frac{\text{No: of elements of class } i}{\text{Total no: of elements in } S}$$

Now

Consider the worst possible case.

It will be $p_S(1) = p_S(2) = \dots = p_S(K) = 1/K$
i.e. equal probability to all classes
it is the most impure set

We want our set to be as far as possible from this set.

Now

let the worst possible set with ' K ' classes in set Q be $p_Q(1) = p_Q(2) = \dots = p_Q(K) = 1/K$

let our set be P with ' K ' classes

and probabilities = $p_P(1), p_P(2), \dots, p_P(K)$

Now

consider a function

$$f(P, Q) = \sum_{i=1}^K P_i \log \frac{P_i}{Q_i}$$

$$f(P, Q) = \sum_{i=1}^K P_p(i) \log \left(\frac{P_p(i)}{P_q(i)} \right)$$

\Rightarrow Let us not go into the details of why this particular function was chosen

\Rightarrow Let us understand what this function tells us

If $P_p(i)$ are equal to $P_q(i)$ then $f(P, Q) = 0$

So the more different $P_p(i)$ are from $P_q(i)$ the larger the function will be.

\Rightarrow Also it can mathematically shown that $f(P, Q) \geq 0$ (always) but let us not go too deep.

\Rightarrow Now our aim is to maximize $f(P, Q)$

$$\Rightarrow \max \left(\sum_{i=1}^K P_p(i) \log \left(\frac{P_p(i)}{P_q(i)} \right) \right)$$

$$\Rightarrow \max \left(\sum_{i=1}^K P_p(i) \log P_p(i) - P_p(i) \log (P_q(i)) \right)$$

$$\Rightarrow \max \left(\sum_{i=1}^K P_p(i) \log P_p(i) - \log \sum_{i=1}^K P_p(i) \log (1/K) \right)$$

$$\Rightarrow \max \left(\sum_{i=1}^K P_p(i) \log P_p(i) - \underbrace{\log (1/K)}_{\text{constant}} \right)$$

$$\Rightarrow \max \left(\sum_{i=1}^K P_p(i) \log P_p(i) \right)$$

$$\Rightarrow \min \left(- \sum_{i=1}^k P_p(i) \log(P_p(i)) \right)$$

Let for set P

$$E(P) = - \sum_{i=1}^k P_p(i) \log(P_p(i))$$

└

This is our entropy as we want to minimize it.

Now Information Gain (IG)

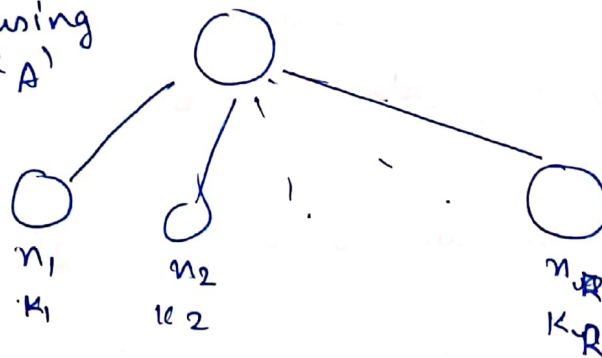
let



node with
n samples
and k diff types
of outcomes

we can split it in many ways depending on our features of our dataset.

let us split using feature 'A'



let it split
into
R nodes

Now $E(n) \rightarrow$ entropy of parent node

$E(n_1) \dots E(n_R) \rightarrow$ entropy of child nodes

Now

IG of this split due to feature 'A' is given as

$$IG(A) = E(n) - \sum_{i=1}^R \underbrace{\frac{n_i}{n}}_{\text{weight of given child node}} E(n_i)$$

weight of given child node.

Now

if gain is high

it means ~~entropy of~~ weighted average of
entropies of child nodes
is much ~~low~~ lower than parent node

So we choose feature with max gain and split
it.

Note \rightarrow Gain ^{can} be +ve, -ve or zero also.