

AWS Deployment  
↓ code base

Github  
Repository

↑  
changes

code pipelines

Pipeline

||  
Automatically  
Deploy

AWS

elastic  
Beanstalk

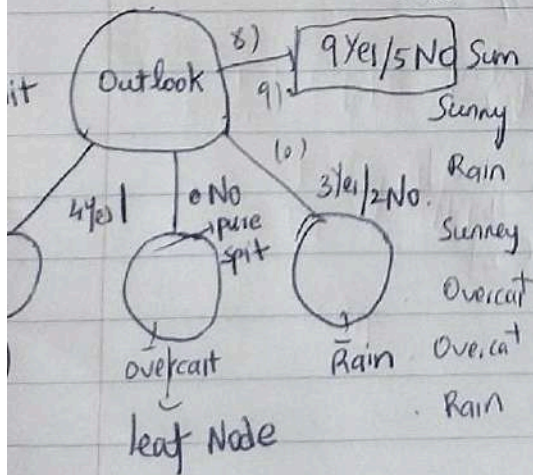
=> Linux machine

↑  
Configuration

## Decision Tree

### Dataset

Day	Outlook	Temperature	Humidity	Wind	Play Tennis
1)	Sunny	Hot	High	Weak	No
2)	Sunny	Hot	High	Strong	No
3)	Overcast	hot	High	Weak	Yes
4)	Rain	mild	High	Weak	Yes
5)	Rain				Yes
6)	Rain				No
7)	Over.				Yes
8)					No





① Purity split check - Pure split (or) Impure split

$\left\{ \begin{array}{l} \text{Entropy} \\ \text{Gini Impurity} \end{array} \right\}$  measure of purity

② What feature you need to select to start the split - Information Gain

① Purity Check

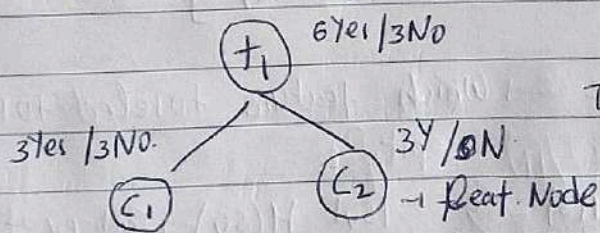
Binary Classification

① Entropy

$$H(S) = -P_+ \log_2 P_+ - P_- \log_2 P_-$$

$P_+$  = probability of positive category

$P_-$  = Probability of negative category



$$H(C_1) = -P_+ \log_2 P_+ - P_- \log_2 P_-$$

$$= -\frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6}$$

$H(C_1) = 1$  - This is impure

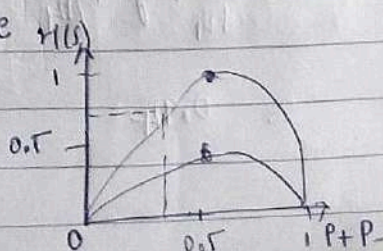
② Gini Impurity

$$G.I = 1 - \sum_{i=1}^h (p_i)^2$$

$P_+ = \frac{3}{6} = \frac{1}{2}$  (Total positive count / Total count)  
 $P_- = \frac{3}{6} = \frac{1}{2}$  (Total negative count / Total count)

$$H(C_2) = -\frac{3}{3} \log_2 \frac{3}{3} - \frac{0}{3} \log_2 \frac{0}{3}$$

$H(C_2) = 0$  pure split





② Gini Impurity  $\rightarrow$  max  $\rightarrow 0.5$

3% low

$$G.I = 1 - \sum_{i=1}^n (p_i)^2$$

$$G.I(C_2) = 1 - \left[ \left(\frac{3}{3}\right)^2 + \left(\frac{0}{3}\right)^2 \right]$$

$$Gini(C_1) = 1 - \left[ (p_1)^2 + (p_2)^2 \right]$$

$$= 1 - 1$$

$C_2 = 0$   $\rightarrow$  Purity split

$$= 1 - \left[ \left(\frac{1}{2}\right)^2 + \left(\frac{1}{2}\right)^2 \right]$$

$$= 1 - \frac{1}{2} = 0.5 = \text{Impure split}$$

Multiclass classification Problem  $\div$  3 categories in O/P

$$H(S) = -p_{C_1} \log_2 p_{C_1} - p_{C_2} \log_2 p_{C_2} - p_{C_3} \log_2 p_{C_3}$$

$$G.I = 1 - \left[ (p_{C_1})^2 + (p_{C_2})^2 + (p_{C_3})^2 \right]$$

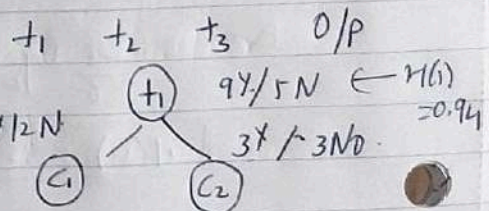
② Information Gain  $\rightarrow$  Which feature is selected to split the split?

$$\text{Gain}(S, t_1) = H(S) - \sum_{v \in \text{val}} \frac{|S_v|}{|S|} H(S_v) \rightarrow \text{Entropy of child node}$$

$$H(S) = -p + \log_2 p + -p \log_2 p -$$

$$= -\frac{9}{14} \log \frac{9}{14} - \frac{5}{14} \log \left(\frac{5}{14}\right) \quad 6 \times 12 \text{ N}$$

$$\approx 0.94$$





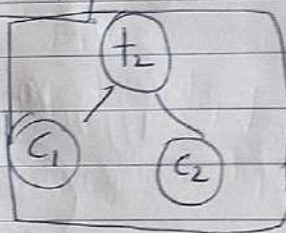
$$H(C_1) = -\frac{6}{8} \log_2 \left( \frac{6}{8} \right) - \frac{2}{8} \log_2 \left( \frac{2}{8} \right) = 0.81$$

$$H(C_2) = -\frac{3}{6} \log_2 \left( \frac{3}{6} \right) - \frac{3}{6} \log_2 \left( \frac{3}{6} \right) = 1$$

$$\text{Gain}(S, f_1) = H(S) - \sum_{v \in \text{val}} \frac{|S_v|}{|S|} H(S_v) \quad \text{Entropy of categories}$$

$$= 0.94 - \left[ \underbrace{\frac{8}{14}}_{\text{Total count}} \times \underbrace{0.81}_{H(C_1)} + \underbrace{\frac{6}{14}}_{\text{Total count}} \times \underbrace{1}_{H(C_2)} \right]$$

$$\boxed{\text{Gain}(S, f_1) = 0.049}$$



$$\Rightarrow \text{Information loss} = 0.051$$

$$\boxed{\text{Gain}(S, t_2) = 0.051} \Rightarrow \boxed{\text{Gain}(S, t_1) = 0.049}$$

We need to start splitting using  $f_2$  feature

~~Entropy~~



## Entropy vs Gini Impurity

When dataset is small  $\rightarrow$  Entropy [log formula]

When dataset is huge  $\rightarrow$  Gini Impurity [simple math]

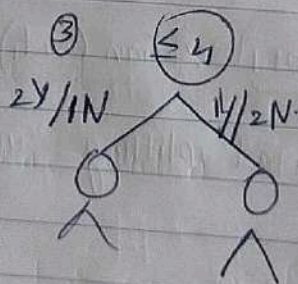
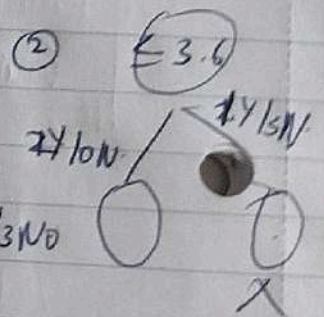
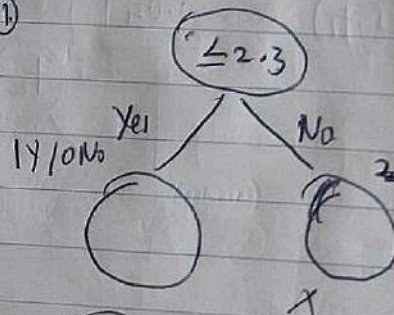
What if my feature is continuous

$t_1$	O/p
2.3	Yes
3.6	Yes
4	No
5.2	No
6.7	Yes
7.8	No

① Sort the feature  $t_1$

① Threshold = 2.3

①



Time Complexity is

huge ~~then~~ when dataset

is huge from sklearn.tree import DecisionTreeClassifier.

classifier = DecisionTreeClassifier()



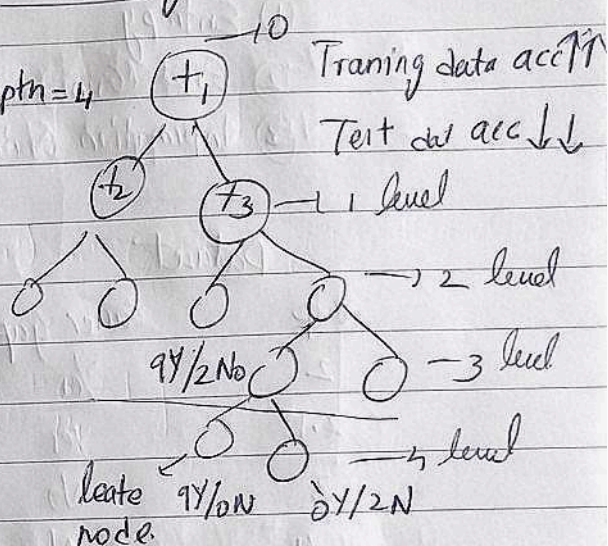
## Decision Tree post pruning and pre pruning [Reduce overfitting]

Training Data

$t_1$   $t_2$   $t_3$  O/p

Generalized model  
Reduce overfitting

Max depth = 4



### ① Post Pruning

- ① Construct the entire decision Tree to complete leaf node.
- ② Pruning the decision
- ③ for suitable for smaller Dataset

### ② Pre pruning

- ① No Hyperparameter Tuning to select Best parameters  
(gridcv, Random Rcv)



## Decision Tree Classifier

- ① Entropy
- ② G-I (Gini)
- ③ Information Gain

## Decision Tree Regressor

- ① Variance Reduction
- ② Variance

### Dataset

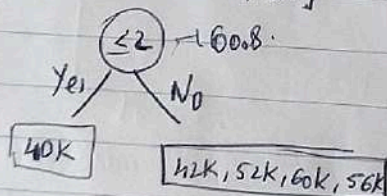
exp	career gap	Salary
2	Yes	40k
2.5	Yes	42k
3	No	52k
4	No	60k
4.5	Yes	56k

$$\bar{y} = 50k$$

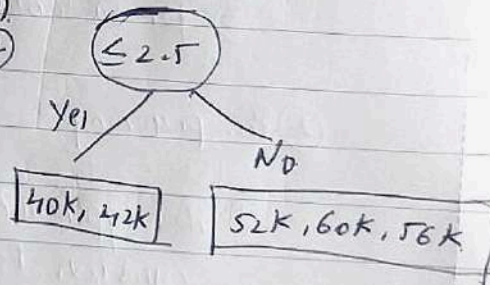
← avg. salary (mean)

Variance Reduction

① [40k, 42k, 52k, 60k, 56k]



②



Final aim - Variance reduction

$$\text{Variance error} = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2 \quad [\text{mean squared error}]$$

$$\begin{aligned} \text{Variance of Root} &= \frac{1}{5} \left[ (40-50)^2 + (42-50)^2 + (52-50)^2 + (60-50)^2 + (56-50)^2 \right] \\ &= \frac{1}{5} [100 + 64 + 4 + 100 + 36] \\ &= 60.8 \end{aligned}$$



$$\text{Variance of left Chi} = \frac{1}{1} [(40-50)^2]$$

$$= 100 //$$

$$\text{Variance of right} = \frac{1}{4} [(42-50)^2 + (52-50)^2 + (60-50)^2 + (50-50)^2]$$

$$= 51 //$$

$$\text{Variance of Reduction} = \text{Var}(\text{Root}) - \sum w_i \text{Var}(\text{Child})$$

$$= 60.8 - \left[ \frac{1}{5} * 100 + \frac{4}{5} * 51 \right]$$

let side have only (one) by total element

$$= 0.$$

$$\text{Variance Reduction} = 0$$

Select the greatest Variance Reduction