

$$\|x_1 - x_2\| = \left(\sum_{i=1}^d (x_{1i} - x_{2i})^2 \right)^{1/2}$$

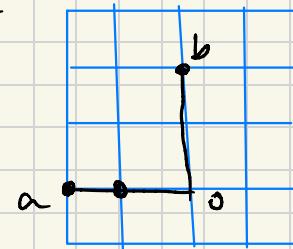
↓
L₂ norm of a vector.

$\|x_1\|$ = Euclidian distance of x_1 from origin $\left(\sum_{i=1}^d x_{ii}^2 \right)^{1/2}$

② Manhattan Distance :

$$\sum_{i=1}^d |x_{1i} - x_{2i}|$$

↑ absolute value



$\|x_1 - x_2\|_1$, L₁ norm of a vector.

$$\begin{aligned} \text{dist}(a, b) &= \text{dist}(a, 0) + \\ &\quad \text{dist}(0, b) \end{aligned}$$

Generalization to L₁ & L₂ norm —

③ L_p norm / Minkowski dist :

$$\|x_1 - x_2\|_p = \left(\sum_{i=1}^d |x_{1i} - x_{2i}|^p \right)^{1/p}$$

if $p=2$ then Minkowski dist = Euclidian dist

if $p=1$ then Minkowski dist = Manhattan dist.

4

Hamming Distance:

XOR between two Boolean vectors.

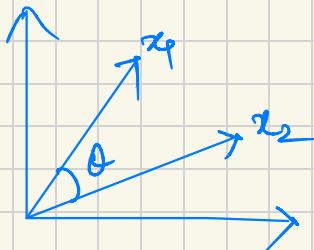
Applicable to strings also,

$$x_1 = \text{a b } \boxed{\text{c}} \text{ a d e f } \quad \text{h i k}$$
$$x_2 = \text{a c } \boxed{\text{b}} \text{ a d e g } \quad \text{h i k}$$

$$\text{hamming dist}(x_1, x_2) = 4.$$

5

Cosine-Similarity:-



$$\cos(\theta) = \frac{x_1 \cdot x_2}{\|x_1\|_2 \|x_2\|_2}$$

if x_1, x_2 is unit vector
then $\cos(\theta) = x_1 \cdot x_2$

6

Edit Distance:-

- The minimum edit distance between two strings.
- 3 operations allowed (Insert, Delete, substitute)
- Needed to transform one into the other.

Rx:

INTE*NTION
|| || || | | |
*EXECUTION

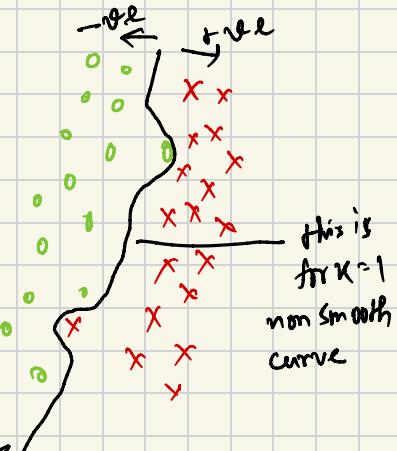
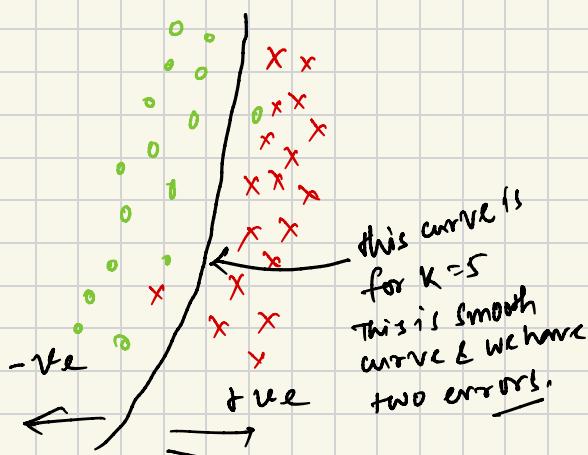
d s s i s

So, distance between these two is 5.

How to decide K?

[K] in KNN is hyperparameter.

geometric intuition if $K=1$, $K=5$



Such curves are called decision surfaces.

If K is too large, let $K=n$

then if $n_1 > n_2$ [$n_1 + n_2 = n$]

In this case, any query point y_m take, it will be +ve. because you have to consider all n points & among $n_1 < n_2$, $n_1 > n_2$. so new print will be assigned to n_1 side.

Overfit

- It tries to cover noisy points also
- non-smooth
- no mistakes

$$K=1$$

Underfit

- Imperfect
- It leaves about majority prints. every points belongs to majority class.
- Imperfect

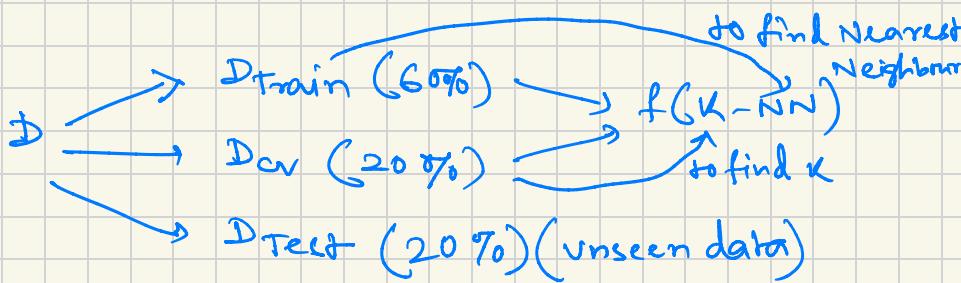
$$K=n$$

Well fit

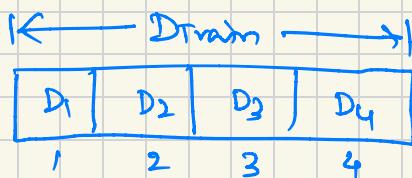
- less prone to noise
- robust
- $K=5$
- smooth.

How to determine K?

① we will devide the whole dataset into 3 parts.



②

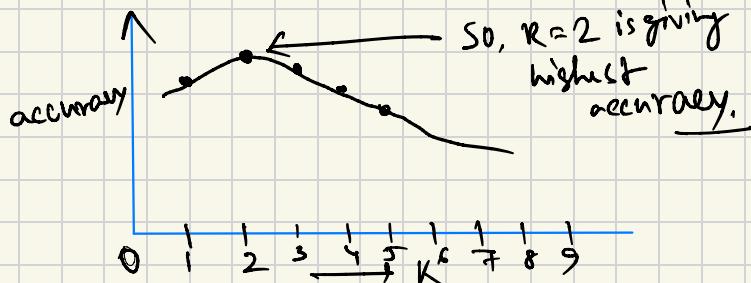


→ randomly divide into 4 parts with equal size. (4-Fold)

③

	Train _m	D _{cv}	Accuracy
K=1	D ₁ D ₂ D ₃	D ₄	a ₁
K=1	D ₁ D ₂ D ₄	D ₃	a ₂
K=1	D ₁ D ₃ D ₄	D ₂	a ₃
K=1	D ₂ D ₃ D ₄	D ₁	a ₄
K=2			a ₁
K=2			a ₂
K=2			a ₃
K=2			a ₄

Considering 4-fold cross-validation we need to plot average accuracy.



Test/Evaluation Time & Space complexity:

Input: $D_{\text{train}}, k, x_q \in \mathbb{R}^d$

Output: y_q \downarrow k is generally small, its 5 or 10.

$KNNpts = []$

for each x_i in D_{Train} :

Compute $d(x_i, x_q) \rightarrow d_i$

Keep smallest k distance (x_i, y_i, d_i)

Count-prs = 0 ; Count-neg = 0

$\rightarrow n$ points, d -dim each
so, time complexity
 $= O(nd)$

\rightarrow considering
 k is very small
it will be $O(1)$

for each x_i in $KNNpts$:

if y_i is +ve :

Count-prs + 1

else Count-neg + 1

if Count-prs > Count-neg:

return $y_q = 1 (+ve)$

else:
return $y_q = 0 (-ve)$

Time Complexity
 $= O(nd)$

Space Complexity

$\approx O(nd)$ to keep
training dataset

** Limitations of KNN:

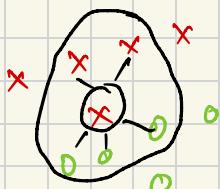
It takes huge time & space. To reduce it we have some techniques

→ KD-tree

→ LSH (Locality sensitive Hashing)

Weighted KNN:

Instead of majority votes, sometimes we can go with shortest distance.



Let $K=5$
for given x_q ,

The way to assign
weights is

$$w_i = \frac{1}{d_i}$$

$d_i \uparrow \rightarrow w_i \downarrow$
 $d_i \downarrow \rightarrow w_i \uparrow$

$$\begin{aligned} \text{Sum} &= 1.5 \\ (1+0.5) & \end{aligned}$$

$$\begin{aligned} \text{Sum} &= 1.75 \\ (1+\frac{1}{2}+\frac{1}{4}) & \end{aligned}$$

Here, $1.5 >> 1.75$ so, it will be assigned
to -ve class.

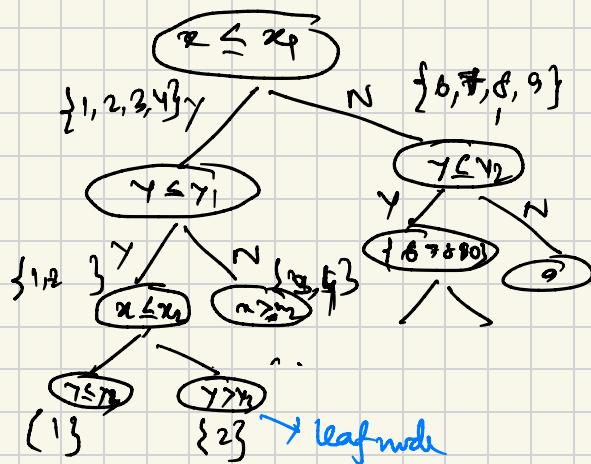
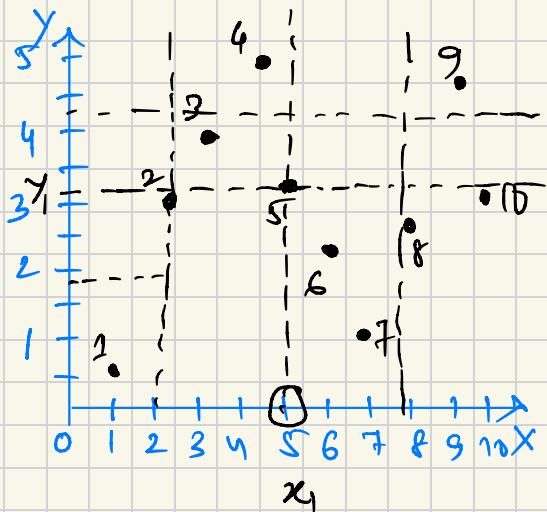
These points are
closest to x_q ,
but if we consider
majority votes,
it will be assigned
to +ve. Now
need to consider
distance also?
Why only majority
votes?
To eliminate this
we can go with
Weighted KNN.

How to implement kd-tree using BST?

BST is \rightarrow 1D [Scalar]

Extending BST to K-dimensional tree
is called Kd tree.

Step 1: Pick X-axis, project points into X-axis, Compute median, split data into median.



You alternate each axis's one after another till leaf nodes.

so, RD-tree boxes the space using axis parallel lines/ planes into rectangle(2D)/ cuboids(3D)/ hypercubrid(n D).

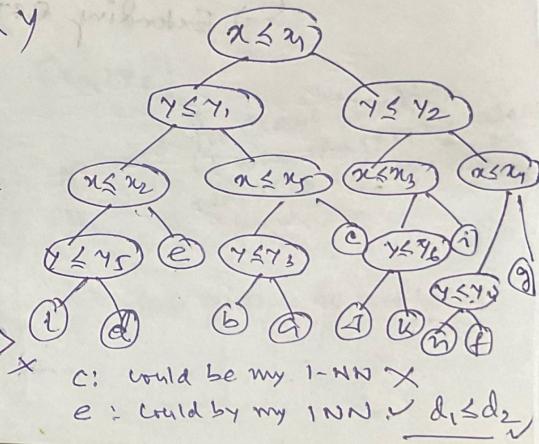
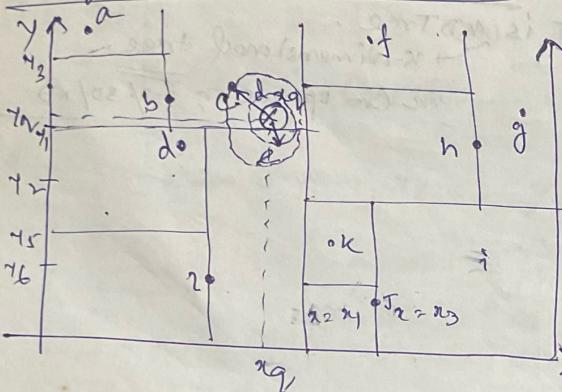
In 2D \rightarrow axis parallel lines \rightarrow rectangles

In 3D \rightarrow axis parallel planes \rightarrow cuboids.

In n D \rightarrow axis parallel hyperplane \rightarrow hypercubrid.

using KD-tree to find INN:

→ using kd tree to find INN:



of comparison to find 1-NN

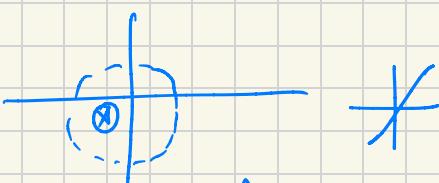
best case: $O(\log n)$

worst case: $O(n)$

If it is K-NN, the best case $O(K \log n)$

worst case $O(K * n)$.

Limitations of KD-tree:



for 2D it need 4 regions to check

When d is not small the time complexity $O(2^d \log n)$

Decision Tree:-

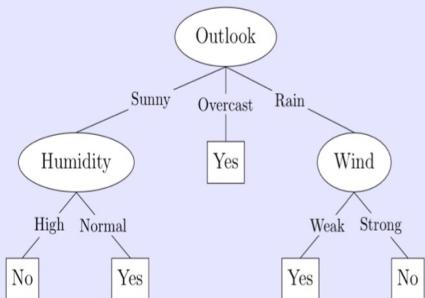
Day	Outlook	Temp	Humidity	Wind	Tennis?
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	No
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No

= {D}



given a new query you just need
to pass your data to DT to
find your output.

→ DT



Algorithm 1 DECISIONTREETRAIN(*data*, *remaining features*)

```
1: guess  $\leftarrow$  most frequent answer in data // default answer for this data
2: if the labels in data are unambiguous then (only single class in the leaf node)
3:   return LEAF(guess) // base case: no need to split further
4: else if remaining features is empty then (if no features left to split then stop)
5:   return LEAF(guess) // base case: cannot split further
6: else // we need to query more features
7:   { for all f  $\in$  remaining features do
8:     NO  $\leftarrow$  the subset of data on which f=no
9:     YES  $\leftarrow$  the subset of data on which f=yes
10:    score[f]  $\leftarrow$  # of majority vote answers in NO
11:      + # of majority vote answers in YES
12:    } // the accuracy we would get if we only queried on f
13:   end for  $\rightarrow$  Storing score of a feature (Information gain)
14:   f  $\leftarrow$  the feature with maximal score(f)
15:   NO  $\leftarrow$  the subset of data on which f=no
16:   YES  $\leftarrow$  the subset of data on which f=yes
17:   left  $\leftarrow$  DECISIONTREETRAIN(NO, remaining features \ {f})
18:   right  $\leftarrow$  DECISIONTREETRAIN(YES, remaining features \ {f})
19:   return NODE(f, left, right)
20: end if
```

How to calculate these scores?

Entropy:

Let's say $y = y_1, y_2, \dots, y_K$ is random variable.

$$H(y) = - \sum_{i=1}^K p(y_i) \log_b(p(y_i))$$

for categorical values

$$H(y) = - \int p(y_i) \log_b(p(y_i)) dy$$

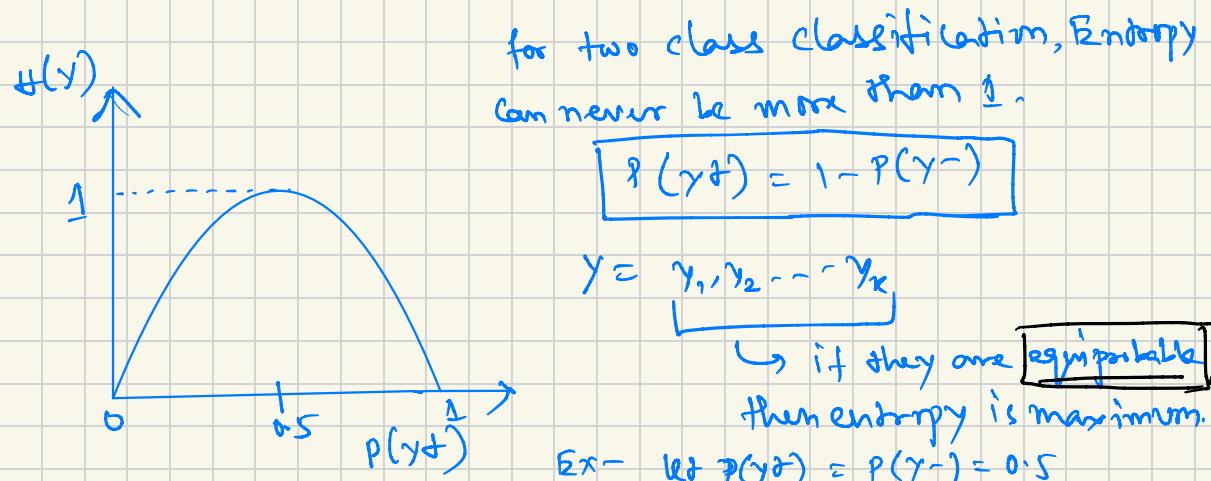
for real values

Here $b = 2$

$$b = 2.718$$

$$\log_2 = \lg$$

$$\log_e = \ln$$



$$Y = Y_1, Y_2, \dots, Y_k$$

↪ if they are equiprobable then entropy is maximum.

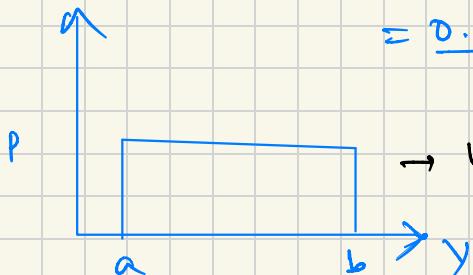
$$\text{Ex- let } P(Y+) = P(Y-) = 0.5$$

$$\begin{aligned} H(Y) &= -P(Y+) \log(P(Y+)) - P(Y-) \log(P(Y-)) \\ &= -0.5 \log 0.5 - 0.5 \log 0.5 \\ &= -(\log 0.5) = -(\log 1 - \log 2) \\ &= \underline{1.} \quad \checkmark \end{aligned}$$

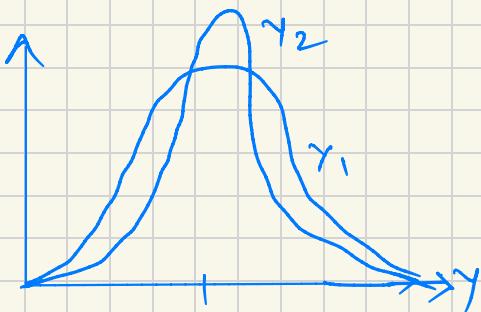
If any one of random variable is most probable & others are 0 then, → entropy is minimal.

$$\text{let } P(Y+) = 1 \quad P(Y-) = 0$$

$$\begin{aligned} H(Y) &= -P(Y+) \log(P(Y+)) - P(Y-) \log(P(Y-)) \\ &= -1 \times 0 - 0 \times \log 0 \\ &= \underline{0.} \end{aligned}$$



→ Uniform distribution, so, if y_i v y_j is equiprobable then it has maximum Entropy.



$\gamma_1 \rightarrow$ less peaked
 $\gamma_2 \rightarrow$ very peaked

So, γ_1 has more data variance
 whereas γ_2 has low data variance

So, $H(\gamma_1) > H(\gamma_2)$

entropy of γ_1 always greater than
 γ_2 .

Partition of Data:

initial Entropy

$$H_D(y) = \frac{8}{14} \times \log\left(\frac{8}{14}\right) - \frac{6}{14} \times \log\left(\frac{6}{14}\right)$$

$$= -0.57 \times (-0.8074) - 0.42 \times (-1.2224)$$

$$= 1.46 + .52$$

$$= .98$$

Let,
 $D \rightarrow D_1$ (5 rows) $H_{D_1}(y) = 0.97$
 $D \rightarrow D_2$ (4 rows) $H_{D_2}(y) = 0$
 $D \rightarrow D_3$ (5 rows) $H_{D_3}(y) = 0.97$

Information gain (IG):

$$\left(\frac{-0.97 \times 5}{14} + 0 \times \frac{4}{14} + 0.97 \times \frac{5}{14} \right)$$

weight entropy after
 partitioning dataset

entropy before
 partitioning

$$IG(y, \text{var}) = H_D(y) - \left[\sum_{i=1}^k \frac{|D_i|}{|D|} * H_{D_i}(y) \right]$$

Gini Impurity % (I_G)

~ similar to entropy

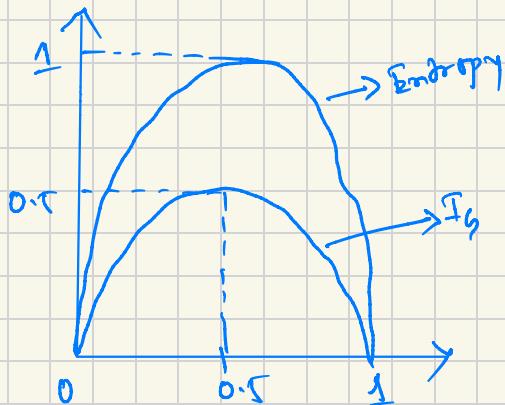
$$I_g(y) \neq I_G(y)$$

for random variable $Y = Y_1, Y_2, \dots, Y_K$

$$I_G(y) = 1 - \sum_{i=1}^k (P(Y_i))^2$$

Case 1: $P(Y+) = 0.5$ $I_g(y) = 1 - [(0.5)^2 + (0.5)^2]$
 $P(Y-) = 0.5$ $= 1 - [0.25 + 0.25]$
 $= 0.5$

Case 2: $P(Y+) = 1$ $I_g(y) = 1 - [1^2 + 0^2]$
 $P(Y-) = 0$ $= 1 - 1 = 0$



max value for Gini Impurity
is 0.5.

$H(y)$ is computationally intensive due to log calculation
whereas Computing Square has less computational cost compare to log computation.

So, I_g is preferable.

Construct a DT

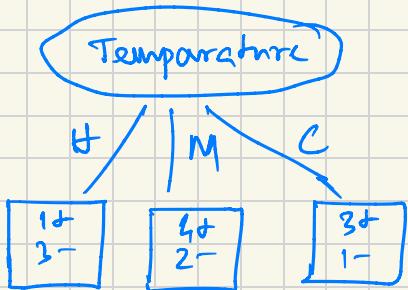
Step 1: $D = \{y^+, y^-\} = \{8, 6\}$ & $H(Y) = .98$

Day	Outlook	Temp	Humidity	Wind	Tennis?
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	No
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No

Step 2:

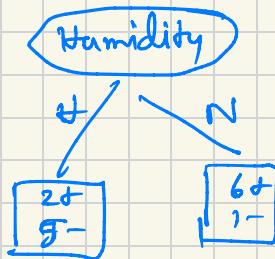
$$\begin{aligned}
 & H(Y) = -\frac{5}{14} \left(\frac{2}{5} \log\left(\frac{2}{5}\right) + \frac{3}{5} \log\left(\frac{3}{5}\right) \right) \\
 & - \frac{4}{14} \left(\frac{3}{4} \log\left(\frac{3}{4}\right) + \frac{1}{4} \log\left(\frac{1}{4}\right) \right) \\
 & - \frac{5}{14} \left(\frac{2}{5} \log\left(\frac{2}{5}\right) + \frac{2}{5} \log\left(\frac{2}{5}\right) \right) \Big] - .98 \\
 & = -\frac{5}{2} \left(\frac{2}{5} \log\left(\frac{2}{5}\right) + \frac{3}{5} \log\left(\frac{3}{5}\right) \right) - \frac{4}{14} \left(\frac{3}{4} \log\left(\frac{3}{4}\right) \right. \\
 & \quad \left. + \frac{1}{4} \log\left(\frac{1}{4}\right) \right) \Big] - .98 \\
 & = -\frac{5}{2} (-.5288 - .4422) + \frac{4}{14} (.3112 - .5) \\
 & = (.3914 + .2312) - .98 \\
 & = (-.6231 - .98) \\
 & = \underline{.35}
 \end{aligned}$$

Step 3:



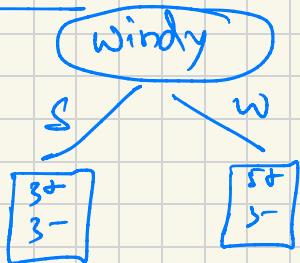
find IG as above

Step 4:



find IG as above

Step 5:



find IG as above

Step 6:

$$IG(y, \text{outlook}) = 35$$

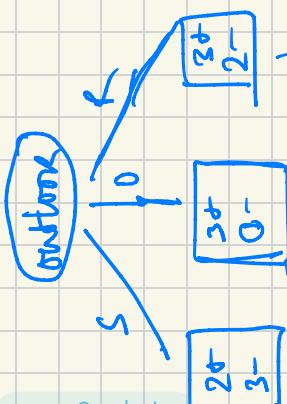
$$IG(y, \text{Temperature}) = -$$

$$IG(y, \text{Humidity}) = -$$

$$IG(y, \text{Windy}) = -$$

We will choose the highest IG as root.

$$IG(y, f) = b_D(y) - \sum_{i=1}^K \frac{|D_i|}{|D|} * H_{D_i}(y)$$



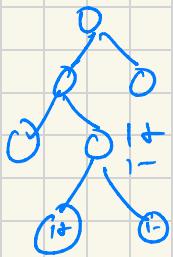
→ no negative points so we can stop here.
This mode is called pure mode.

We are recursively branching each nodes based on 16.

** three points where to stop growing the tree —

① pure node → stop growing the tree.

② Can't grow the tree anymore due to lack of points



In this case no point of growing the node.
Suppose, your dataset has 10k points & you have only one points which follows one path. That point may be noisy. (It will overfit)

③ if we are too deep → stop growing the tree,

As, the depth of the tree ↑ → chance of overfitting ↑

depth is small → underfit.

In DT → depth of the tree is a hyperparameter.

↳ we can use C_V to find it.

Pruning: Remove subtrees for better generalization (decrease variance)

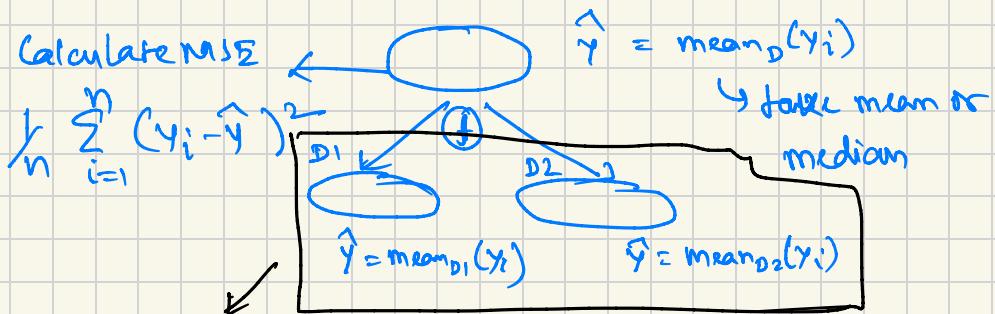
Prepruning: Early stopping

Postpruning: Grow the whole tree then prune subtrees which overfit on the pruning set

Prepruning is faster but postpruning is more accurate.

DT for Regression :

- IG used for classification
- MSE or MAE is used for regression



find weighted mean square

$$w_1 \text{MSE}_{D1} + w_2 \text{MSE}_{D2}$$

Now subtract $\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y})^2 - [w_1 \text{MSE}_{D1} + w_2 \text{MSE}_{D2}]$

whichever split minimizes MSE, should consider for split.

Dimensionality Reduction

volume of a d-dim hyperspace

$$V_d(r) = \frac{r^d \pi^{d/2}}{\Gamma(\frac{d}{2} + 1)}$$

$$V_3(r) = \frac{r^3 \pi^{3/2}}{\Gamma(\frac{5}{2})}$$

$$\Gamma(\frac{5}{2}) = \int_0^{\frac{\pi}{2}} e^{-x} x^{\frac{5}{2}-1} dx$$

$$\Gamma(1 + \frac{3}{2})$$

$$V_3(r) = \frac{r^3 \pi^{3/2}}{\frac{3}{4} \pi^{1/2}}$$

$$= \frac{3}{2} \sqrt{\frac{3}{2}}$$

$$V_3 r = \frac{4}{3} \pi r^3 \quad \checkmark$$

$$= \frac{3}{2} \times \frac{1}{2} \sqrt{\pi}$$

$$= \frac{3}{4} \sqrt{\pi}$$

$$\frac{V_d(\text{spare})}{V_d(\text{cube})} = \frac{r^d \pi^{d/2}}{\frac{(2\pi)^d}{2^d} \Gamma(\frac{d}{2})}$$

Gamma function

$$\Gamma(n) = \int_0^\infty e^{-x} x^{n-1} dx$$

$$\Gamma(1) = 1, \quad \Gamma(\frac{1}{2}) = \sqrt{\pi}$$

$$\Gamma(n+1) = n \Gamma(n)$$

$$\Gamma(n+1) = n!$$

$$\textcircled{a} \quad \Gamma(n) \Gamma(1-n) = \frac{\pi}{\sin n\pi}$$

