

1	2	3	4	5	6
3	5	5	5	5	3

24 30	20 30
----------	----------

Name: Anirban Majumdar	Roll Number: MKS202304
Subject: DMM2	Date:
Course & Year:	Total No. of Pages:

— Begin here —

1) number of transaction $(T) = 10^{10}$

total no of elements (not necessarily unique) $av \leq \frac{10^{10} \cdot 10}{10^4}$
 $= 10^{11}$

set of items has size (10^7)

an item is frequent if it appears in more than

$$\frac{10^{10}}{10000} \text{ transaction} \\ = (10^7) \text{ transaction}$$

$$\text{So, at most } \frac{10^{11}}{10^7} \text{ items can be frequent} \\ = 10^4 \quad \checkmark$$

So total 10^4 items can be frequent

$$|F_1| \leq 10^4$$

Now C_2 will have form $= \{(x_1, x_2) \mid x_1 \in F_1 \wedge x_2 \in F_1\}$

$$\text{So, } |F_2| \leq |C_2|$$

$$\leq \binom{10^4}{2} \\ = \frac{10^4(10^4-1)}{2} \\ = 5 \cdot 10^3(10^4-1)$$

Better estimate by
directly calculating F_2 ,
like F_1

(3)

(4) In random forest classifier,

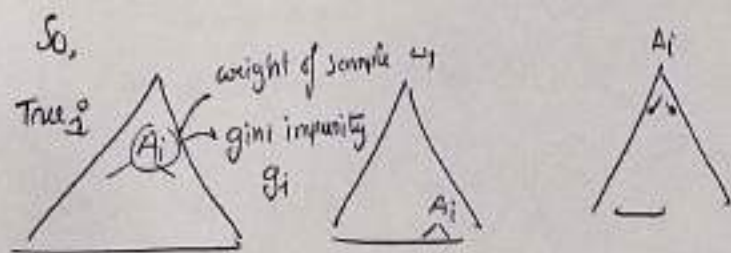
1 set of attributes from which to choose

there will be multiple decision trees and further if m is low,

the correlation and strength will be lower. So the attribute on which the tree splits consecutively changes,

i.e. the attribute (splitter) will come at various heights of

the different trees. (may not come too).



$$\text{So, Gini Impurity}(A_i) = \frac{\sum_j (w_j \cdot g_j)}{\sum_j w_j} \quad \text{if } A_i \text{ appears in True}_j$$

Not only every A_i has a good chance to appear as a splitter but also, its measured at different levels.

Weighted average

Then, they can be ranked based on the Gini Impurity.

But for decision tree, some input features may not occur or

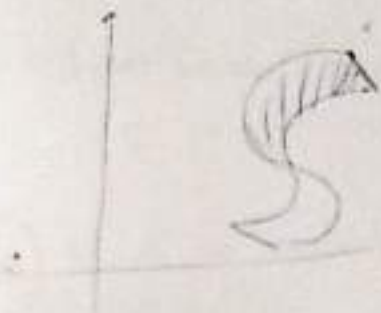
they have only one chance in a path to be appeared and the one with more IG is chosen as the splitter.

A splitter may be good or bad at certain levels, decision tree doesn't handle that case. Some effect for Random forest.

(5)

(3)

6. locally linear embeddings are better alternative for PCA if the data points / are not linear (lie along a manifold)



In the above case, let's say there are m datapoints as x_j and

we take k nearest datapoints of it

$$\text{we say discrepancy of } x = \left(x_i - \sum_{j=1}^k w_{ij} x_j \right)$$

$$= \left(x_i - \sum_{j=1}^m w_{ij} x_j \right)$$

also $w_{ij} = 0$ if x_j not in neighbourhood of x_i

Now we will in general try to minimize the ^{Mean} squared or Mean discrepancy, i.e. the (SSE equivalent) or MSE.

whole dataset:

$$\text{So, } J = \frac{1}{2n} \sum_{i=1}^n \left(x_i - \sum_{j=1}^m w_{ij} x_j \right)^2$$

MSE

Now \hat{W} is the weight matrix s.t.

$$\hat{W} = \underset{W}{\operatorname{argmin}} J$$

The min can be achieved by gradient descent or direct calculation

Once, we get the \hat{W} matrix.

for each unseen point, we will get K nearest neighbours
get the labels for them and then
predict class on based on K nearest neighbours and
weight matrix



Weights \rightarrow Embedding?

(3)

(3)

(A₁)(A₂)

From assuming an attribute ~~is~~ is more significant than other.

if we only use A₁ to predict y, the result will be better than if we have used A₂ as the only attribute

while
let, A₁ has parameter θ_1 , A₂ has θ_2
In that case, we know 1st component of vector x

$$\frac{1}{2} \sum (y_i - \theta_1 x_i^1)^2 \leq \frac{1}{2} \sum_{i=1}^n (y_i - \theta_2 x_i^2)^2$$

(2nd component of vector x)

$$\Rightarrow \sum [(y_i - \theta_1 x_i^1)^2 - (y_i - \theta_2 x_i^2)^2] \leq 0$$

$$\Rightarrow \sum [(2y_i - \theta_1 x_i^1 - \theta_2 x_i^2)(\theta_2 x_i^2 - \theta_1 x_i^1)] \leq 0$$

for this condition to be true,

its not enough to make θ_1 larger than θ_2 as x_i^1 may have higher values compared to x_i^2 .

To illustrate this case consider two tables

(A₁)

A ₀	Price in Rupees	A ₂	If the product good

(A₂)

A ₀	Price in pairs	A ₂	If the product good

lets say price in rupees is the most important so, $\theta_1 \geq \theta_2$ and $\theta_1 \geq \theta_0$
but, PTO

The same data can be used with ^{inputs} changed to μ and σ
and obviously, the prediction will be same

So, $\sigma_1' = \frac{\sigma_1}{100}$ which not necessarily mean $\sigma_1' > \sigma_2$
and $\sigma_1' > \sigma_0$

▣ The second interpretation could be
an attribute is important if ^{small change in} it will make a good change
in prediction values.

By ^{this} definition, more important features will have higher value of σ .

Considering both the cases,

whether importance of σ , (magnitude of σ) implies most significant
feature depend on the definition of significance of attribute

OK (5)

(5) Some observation:

let's take two points (x_i, x_j) among the dataset



⊛ distance

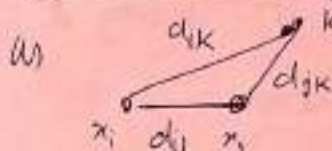
if x_i and x_j are small enough. (wrt the average distance between points)

the two vectors will be close enough in N dimension

Intuition:

	d_i	d_j
1	d_{i1}	$\frac{1}{2} \begin{cases} d_{i1} \pm \delta \\ \leq d_{i1} \end{cases}$
2	d_{i2}	$\frac{1}{2} \begin{cases} d_{i2} \pm \delta \\ \leq d_{i2} \end{cases}$
...
j	d_{ij}	d_{ij}
...
N	d_{iN}	d_{iN}

(good approximation for)
the $\|s\|$ will be d_{ij} itself



$$\delta \leq d_{ij}$$

$$d_{ij} + d_{jk} > d_{ik}$$

as d_{ij} is small,

$$(d_{jk} - d_{ik} < d_{ij})$$

two vectors will be in close vicinity, then distance will be

$$\leq \sqrt{N} (d_{ij})^2$$

$$= d_{ij} \sqrt{N}$$

between x_i and x_j

⊛ if distance is large, then the how large it will be, depending on

1) d_{ij}

2) the structure of cluster

but both the cases, it will be much larger compared to the 1st case

Induction: 1) If the x_i, x_j points are far apart compared to average distance

then, assuming random distribution points again d_{ik} and d_{jk}

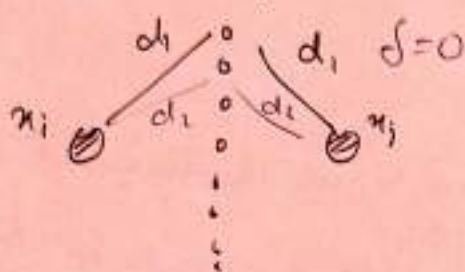
will be ^{a lot} ~~very~~ different, therefore increasing their distance in N dimension

In this case too, $f \leq (d_{ij})$

but \uparrow itself is large.

\Rightarrow

If we consider some specially created distribution,



In this case d_{ij} will be large enough and that will be the only one contributing to the distance of d_i and d_j in N dimension

So, if we cluster this way,

if x_i and x_j are nearby, $\text{dist}(d_i, d_j)$ is len
so d_i and d_j are nearby

if x_i and x_j are large, d_i and d_j are at distance too
though distance depends on the structure
(which might be actually good)

Considering these, for general clusters, clusters formed by column ^{of} D will be ^{an} overall good approximation of X

Class association rules have a ^{attribute} class (attribute) as their target.

Now, if we take a sample decision tree,



Class label $\{C_1, C_2, C_2\}$

Now, we can generate class association rules as follows:

Take a path (P_1 from figure), then

$$(A_1 = a_1^0, A_2 = a_2^0, A_3 = a_3^0) \longrightarrow C_1$$

for all such root to leaves path, we will get such association rules

Now, to generalize a class association rule, we can do the following,

$$(A_1 = a_1^0, A_2 = a_2^0, A_3 = a_3^0) \longrightarrow C_1 \text{ for } 2/2 \text{ samples}$$

$$(A_1 = a_1^0, A_2 = a_2^0, A_3 = a_3^1) \longrightarrow C_2 \text{ for } 1/101 \text{ samples}$$

$$(A_1 = a_1^0, A_2 = a_2^0, A_3 = a_3^2) \longrightarrow C_2 \text{ for } 1/202 \text{ samples}$$

$$\text{So, } (A_1 = a_1^0, A_2 = a_2^0) \longrightarrow C_2 \text{ for } 3/305 \text{ samples}$$

removing A_3 as a splitter
In this case, we have generalized the association rules by clubbing
them into 1.

This is equivalent to pruning the last level (not asking
the last splitting question)

▣ This would be different as, ~~if we club~~ in general ~~we club~~
or prune from bottom (as was the previous case).

what if class association generalization rule has chosen
to remove A_1 instead of A_3 to generalize better
Bottom up pruning won't help, in such case
(Top down pruning/clubbing) might be helpful.

remove A_1 from the class association rule \Rightarrow
change the ^{is a splitter} question to based on attribute (A_2 and the A_3)

on decision tree
So, class association generalization is more exhaustive wrt
usual bottom up pruning.