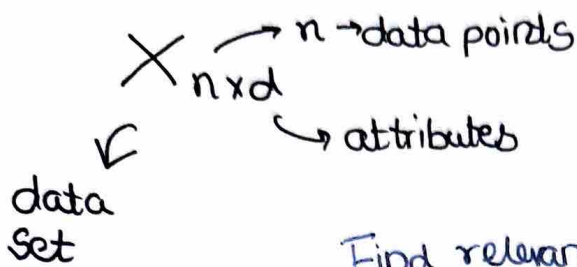


Precision, Recall and F1 Score



Find relevant attributes, which are subset of  $X_{n \times d}$

Why?

Ex:-

medical Dataset

Figure out Vital Signs (out of entire dataset)

→ When smaller no of point suffice

→ use smaller models, faster models

Reasons

→ Interpretability

→ Faster models in Deployment

→ Better Generalization

Eg:- Stock market Data

$10^6$  attributes (noisy classification)

→ 10 most important (better classification)

Very Complex Loss Function

very minimal effect

How to do attribute selection?

① Exhaustive Search (all subsets)  
 $(10^6)^{10} \sim 10^6$

(revise)

$L_1$ -regularization

Lot of learned weights of the model close to 0  
→ those attributes are not important

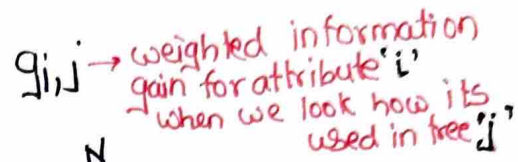
$L_1$ -regularization

$$\text{Loss } L_2 = \text{loss}_{\text{original}} + \lambda \sum w_i^2$$

regularization

Ex:-  
Weights [15, 16, 1, 2, 19]  
↓  $L_1$  (-2)  
[13, 14, 0, 0, 17]

### \* How to use random forests to do feature selection?



$$G_i = \sum_{j=1}^N g_{i,j} \rightarrow$$

$$\hat{G}_{il} = G_{il} / (\sum G_{il}) \quad \text{Importance Scores}$$

$i \rightarrow$  attribute (gain)

normalize the scores

{ attributes which have  
a higher importances } Considered

Random Forest  $\rightarrow$  Collection of Decision Trees.  
other attributes below can be repeated. There will be gain at attributes below too.

- Correlated { (correlation dilutes importance) }  $G_{i1} = \frac{G_i}{\sum_i G_i} \rightarrow$  gain of  $i^{\text{th}}$  attribute  
 { one feature can substitute another }  $\rightarrow$  summation of all gains across all attributes.

example we have :-

example we have :-

attributes

Temperature  
(Celsius)

Temperature  
(Fahrenheit)

→ convey the same thing  
(collectively contribute high importance)

(higher deg polynomial)

correlated

- ↳ numerical values which are bucketized
- ↳ considered important when they aren't

each one of them would be 50% of the tree. (we may miss important attributes)

# \* Select the weakest attribute :- (Greedy Feature Elimination)

Assume for simplicity we have 4 attributes

$X_{n \times 4}$

k-nn algorithm (Classification)

Classifiers

accuracy

C1:  $X / \text{Attribute} = 1$  90%

C2:  $X / \text{Attribute} \# 2$  95%

C3:  $X / \text{Attribute} \# 3$  83%

C4:  $X / \text{Attribute} \# 4$  96%

See their performance on train and test without the corresponding attribute

Evaluate which is the least important

Highest accuracy

C4  $\rightarrow$  without attribute #4



attribute 4  $\rightarrow$  not important  
(we can still reach high accuracy dropping 4)



Drop 4

Ablation Testing

$\rightarrow$  These classifiers are like

C1  $\rightarrow$  without attribute #1

C2  $\rightarrow$  without attribute #2

C3  $\rightarrow$  : : #3

C4  $\rightarrow$  : : #4

(Greedy Feature Elimination)



Requires lot of training



Gives smaller dataset

$$X' = X / 4$$



$X_{\text{prime}}$



$X$  without attribute #4

If 'd' attributes



we need to train

'd' times

$\times$  attributes

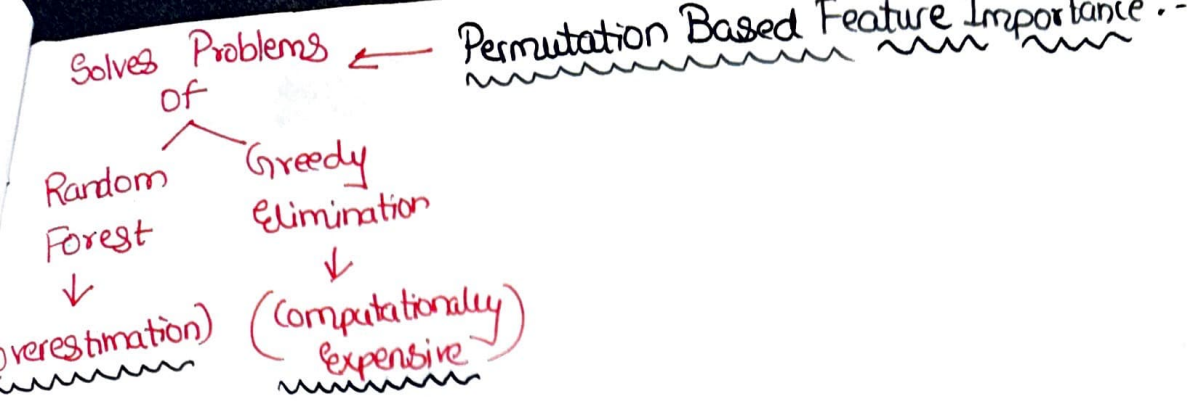
if we want to drop 'm' attributes

$m \cdot d$

m times

$\rightarrow$  downside (Computationally expensive)





Inputs - fitted predictive model  $m$ , tabular dataset (training or validation)  $D$

Compute the reference scores  $s$  of the model  $m$  on data  $D$  ( $n \times d$ )  
 (for instance the accuracy for a classifier or the  $R^2$  for a regressor)

For each feature  $j$  (column of  $D$ ):

- For each repetition  $k$  in  $1, \dots, K$ :
  - Randomly shuffle column  $j$  of dataset  $D$  to generate a corrupted version of the data named  $\tilde{D}_{k,j}$
  - Compute the scores  $s_{k,j}$  of model  $m$  on corrupted data  $\tilde{D}_{k,j}$
- Compute importance  $i_j$  for feature  $f_j$  defined as:

$$i_j = s - \frac{1}{K} \sum_{k=1}^K s_{k,j}$$

$\left. \begin{array}{l} \text{mean of all shuffles (k)} \\ \text{Accuracy} = 85\% \text{ (0.85)} \end{array} \right\} \text{d times / evaluation}$

Shuffled  $k$  times

Explanation

features

①	1.4	②	③
	1.1		
	-0.1		
	2.2		

X

0
2
1
0

Y

Accuracy = 99% (0.99)

1.1
-0.1
2.2
1.4

X'

0
2
1
0

Y

Importance =  $0.99 - 0.85$  (higher accuracy loss)  
 (ot #1) = 0.14 (High Importance)