



RANDOM FORESTS

Nissy Milcia William

nissymilcia.w@niser.ac.in

Under the guidance of Dr. Subhankar Mishra

MACHINE LEARNING CS-460

National Institute of Science Education and Research Bhubaneswar



RANDOM FOREST

A random forest is a machine learning algorithm consisting of a collection of tree structured classifiers where each tree casts a unit vote for the most popular class at some

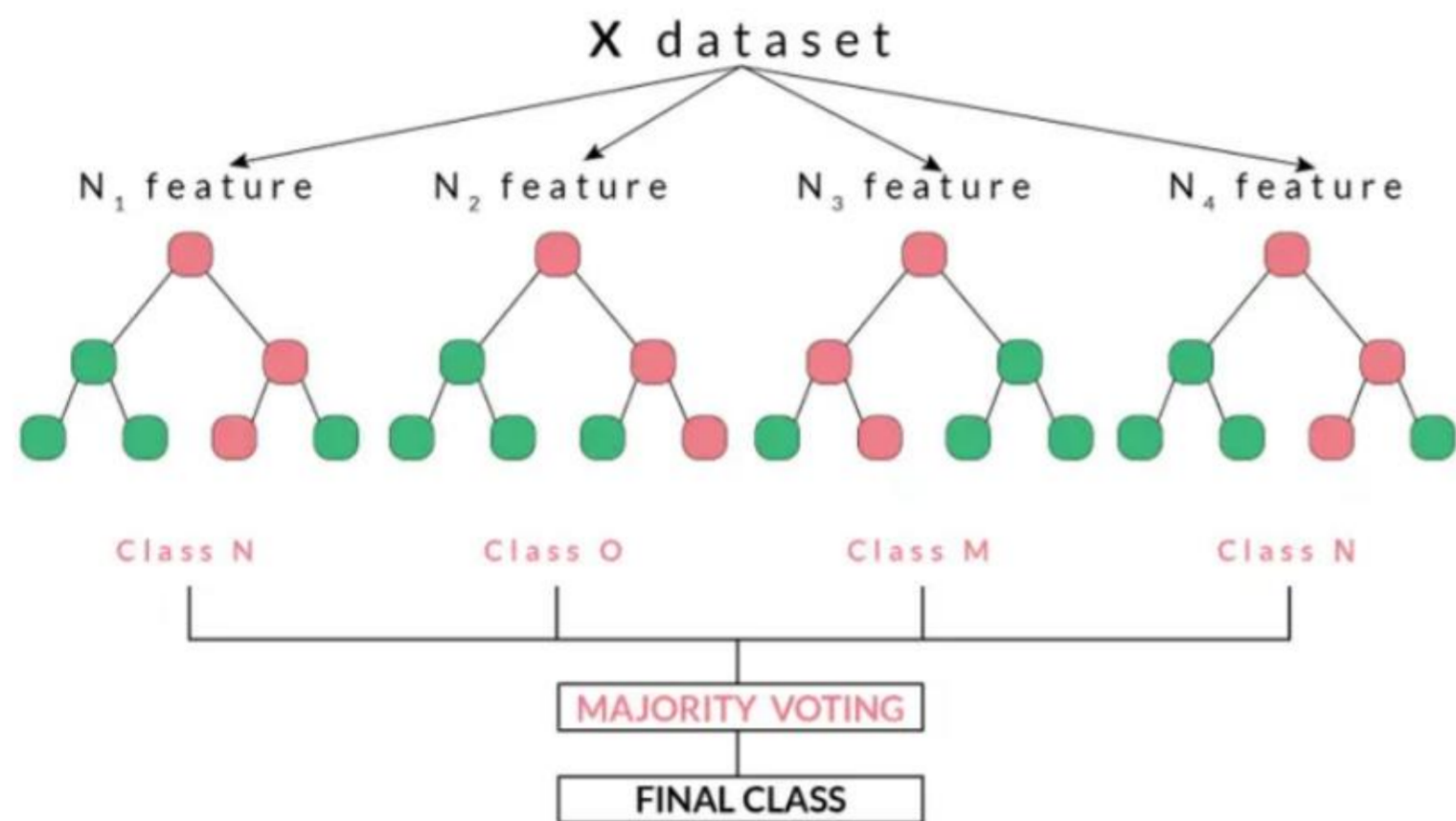
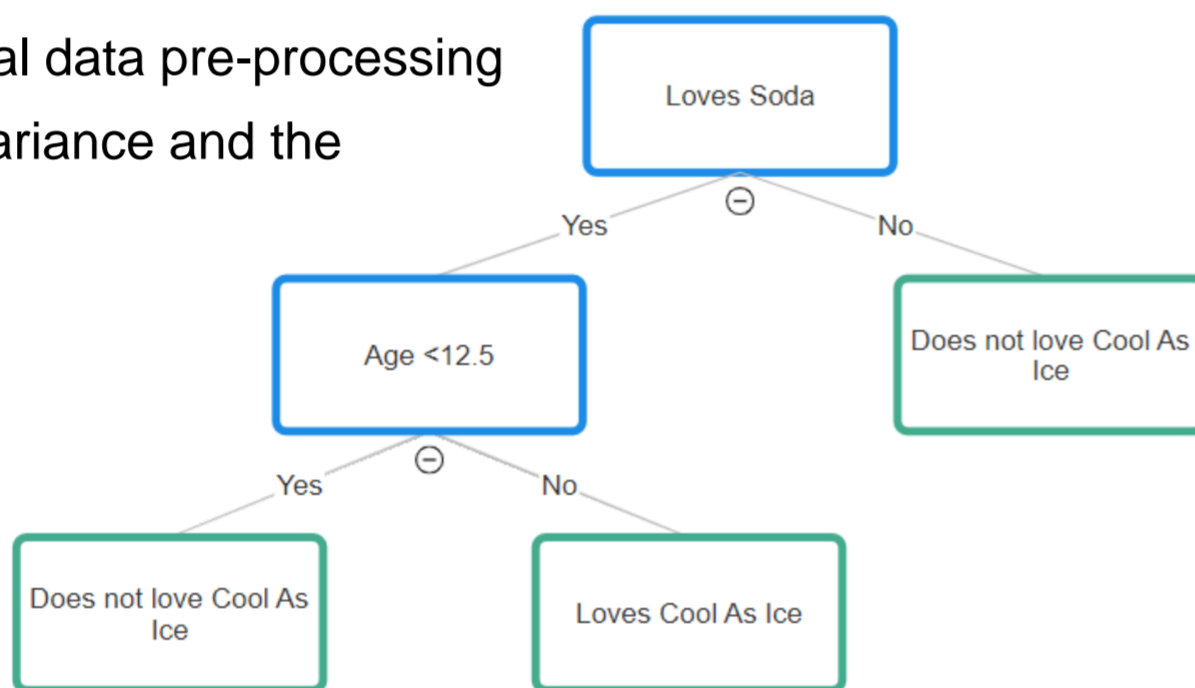


Image from <https://images.app.goo.gl/HzKPL4tH5uvjosnD9>

DECISION TREES

- Flowchart-like tree structures representing the feature, the rules and the result of the algorithm
- Intuitive approach, requires minimal data pre-processing
- Prone to overfitting as a result of variance and the entire structure of the tree can

Loves Popcorn	Loves Soda	Age	Loves Cool As Ice
Yes	Yes	7	No
Yes	No	12	No
No	Yes	18	Yes
No	Yes	35	Yes
Yes	Yes	38	Yes
Yes	No	50	No
No	No	83	No



ENSEMBLE LEARNING: BAGGING

- Group of models (weak learners) work to achieve a final prediction, rather than a single model
- Consists of Bootstrapping, Parallel training and aggregation

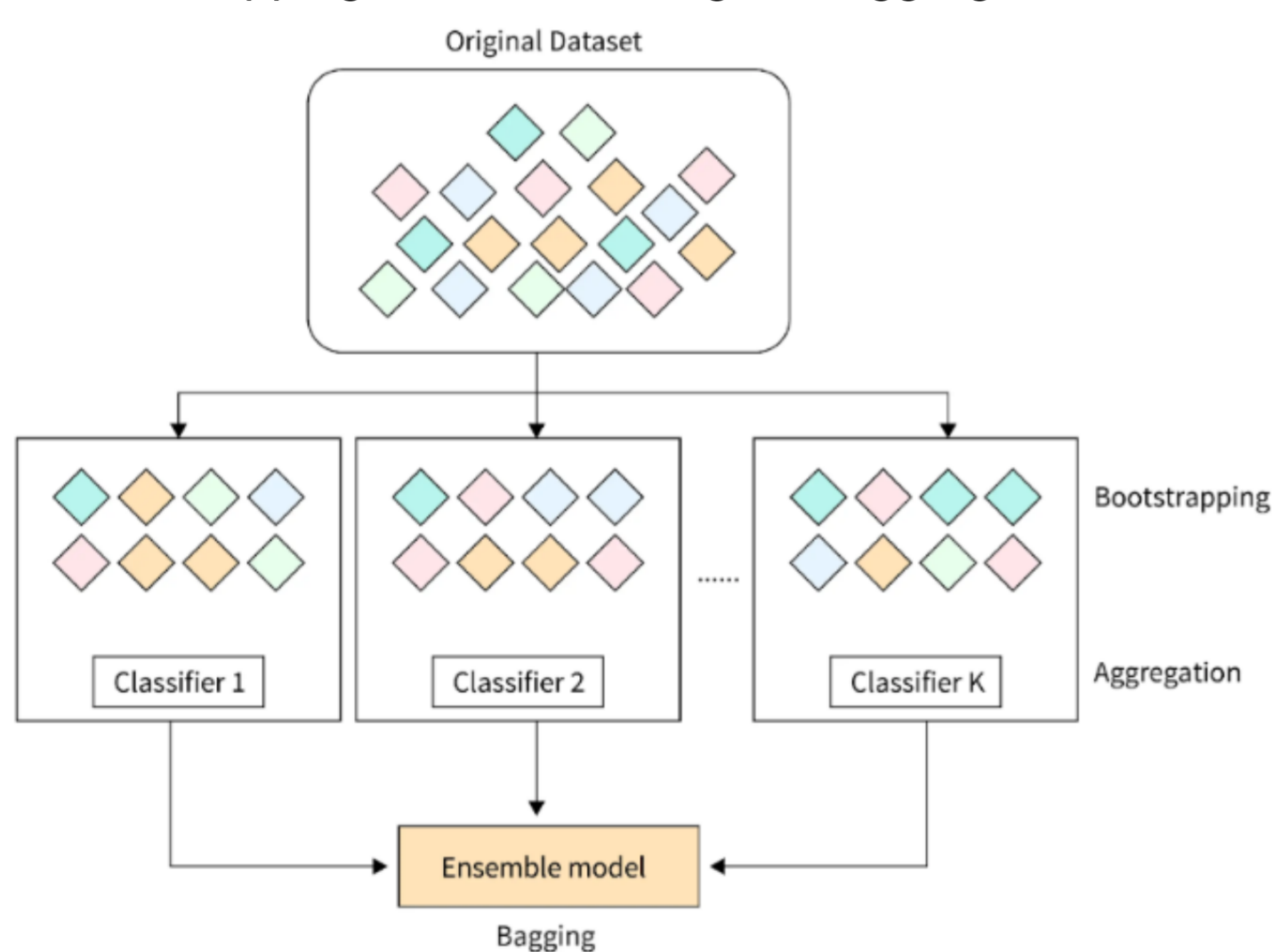


Image from: <https://images.app.goo.gl/PupB1pzREbiGnUZj8>

WHY RANDOM FOREST?

- Ensemble of multiple decision predicts more accurate results, particularly when the individual trees are uncorrelated
- Randomness helps increase tree diversity and is generated by:
 - Bagging:** Ensures that every tree is built from a different subset of the original dataset, can give ongoing estimates of the generalization error of the combined ensemble of trees (out-of-bag estimates)
 - Random Feature selection:** During the selection of the appropriate variable to split on at each node, a fixed number of random variables are considered, rather than the entire collection of variables
- Out-of-bag estimates**

ALGORITHM

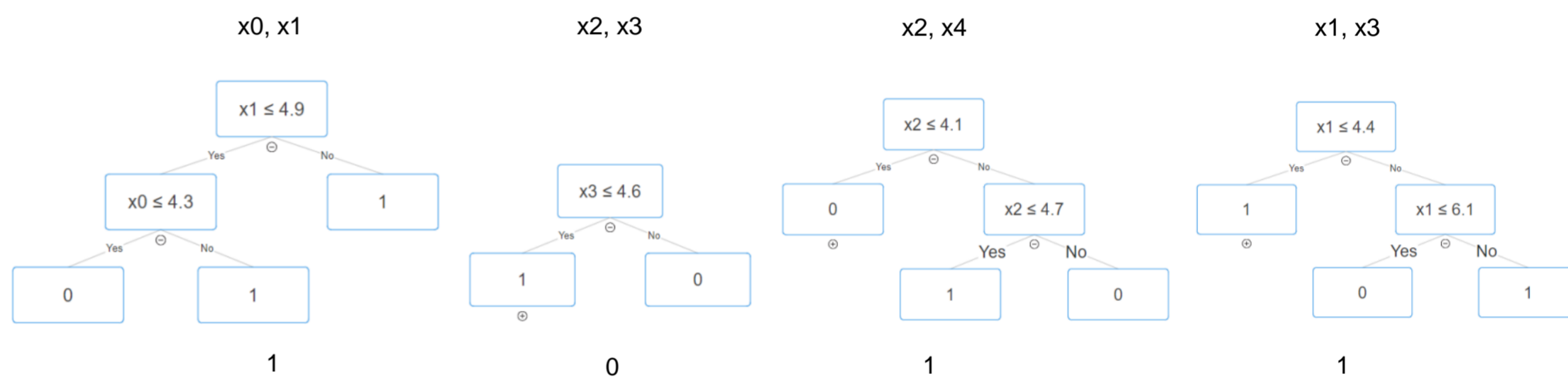
```

Precondition: A training set  $S := (x_1, y_1), \dots, (x_n, y_n)$ , features  $F$ , and number of trees in forest  $B$ .
1 function RANDOMFOREST( $S, F$ )
2    $H \leftarrow \emptyset$ 
3   for  $i \in 1, \dots, B$  do
4      $S^{(i)} \leftarrow$  A bootstrap sample from  $S$ 
5      $h_i \leftarrow$  RANDOMIZEDTREELEARN( $S^{(i)}, F$ )
6      $H \leftarrow H \cup \{h_i\}$ 
7   end for
8   return  $H$ 
9 end function
10 function RANDOMIZEDTREELEARN( $S, F$ )
11   At each node:
12      $f \leftarrow$  very small subset of  $F$ 
13     Split on best feature in  $f$ 
14   return The learned tree
15 end function

```

The predictions from each tree are combined either by a majority vote (for classification tasks) or an average (for regression tasks)

id	x0	x1	x2	x3	x4	y
0	4.3	4.9	4.1	4.7	5.5	0
1	3.9	6.1	5.9	5.5	5.9	0
2	2.7	4.8	4.1	5	5.6	0
3	6.6	4.4	4.5	3.9	5.9	1
4	6.5	2.9	4.7	4.6	6.1	1
5	2.7	6.7	4.2	5.3	4.8	1



APPLICATIONS

- Random forest algorithms are versatile and can be used for classification as well as regression tasks
- Capable of handling large datasets, missing and complex data. They are regularly applied across fields such as Finance, Healthcare, E-commerce, Manufacturing. Here are some specific examples:

A top-down approach to classify enzyme functional classes and sub-classes using random forest

Chetan Kumar* and Alok Choudhary

Using Random Forest Classifier for Particle Identification in the ALICE Experiment

Tomasz Trzeciński¹, Lukasz Graczykowski², and Michał Glinka¹ for the ALICE Collaboration

Random Forest Models To Predict Aqueous Solubility

David S. Palmer, Noel M. O'Boyle,† Robert C. Glen, and John B. O. Mitchell*

Unilever Centre for Molecular Science Informatics, Department of Chemistry, University of Cambridge, Lensfield Road, Cambridge CB2 1EW, United Kingdom

Received May 5, 2006

REFERENCES

- https://youtu.be/_L39rN6gz7Y?si=ZiVx1xZG3k4iRvFz
- <https://youtu.be/RtrBtAKwxcQ?si=t8MGQuXc5a2jvAja>
- <https://www.ibm.com/topics/bagging#:~:text=Bagging%2C%20also%20known%20as%20bootstrap,be%20chosen%20more%20than%20once.>
- <https://www.ibm.com/topics/random-forest>
- <https://youtu.be/sQ870aTKqIM?si=FJYLdc5s3S-5bUqu>
- https://youtu.be/J4Wdy0Wc_xQ?si=IMa6mX41G9cakktc
- <https://www.stat.berkeley.edu/~breiman/randomforest2001.pdf>
- <https://www.geeksforgeeks.org/cart-classification-and-regression-tree-in-machine-learning/>
- <https://link.springer.com/article/10.1023/A:1010933404324>
- <https://pages.cs.wisc.edu/~matthewb/pages/notes/pdf/ensembles/RandomForests.pdf>